

Oilwell Production and Completion Regression Analysis

Tom Wilson

November 24, 2018

Introduction

For centuries, oilwells have been drilled vertically and completed by perforation only. Since the invention of horizontal drilling in the late 20th century, literature has shown that various completion parameters have a significant effect on oil production. Multiple recent sources agree that there is value in leveraging available data as input to statistical models steeped in machine learning.(Fulks et al. 2016), (Holdaway, Laing, and others 2015), (Pankaj et al. 2018), (Subrahmanya et al. 2014) While there is some disagreement about the limitations and applicability of such models, most agree that a general data analysis workflow is applicable. (Fulks et al. 2016),(Temizel et al. 2015),(Groulx et al. 2017), (Pankaj et al. 2018)

Data Quality

One recurring theme within recent literature related to oilwell completion is that of data quality. Lopes and Jorge (2017) takes a statistical approach to dealing with gaps in data by breaking synthesizing gaps in otherwise complete data. Lopes and Jorge (2017) eventually shows that gaps in data did not have a significant effect on model selection or accuracy. On the other hand, Wang et al. (2016) claims that a input uncertainties model selection challenges were primarily caused by missing data and noise. Khodabakhsh, Ari, and Bakir (2017) took yet another approach via real time models for detecting and classifying errors. Wang et al. (2016) suggests Basis Pursuit Denoising (BPDN) as a solution to data quality issues. Other sources did not comment on dealing with data quality or noise.

General Workflow

Regardless of the target variable or underlying dataset, there is wide-spread agreement that data analysis workflows are generally applicable and value-added. Quoting from Holdaway, Laing, and others (2015): “Data-driven workflows, models, and analysis can address a diverse array of business problems in the oil and gas industry.” Groulx et al. (2017) found that their approach applied equally well to all basins(a geologic formation) and plays(oil bearing zone) available. Groulx et al. (2017) showed that, generally, the number of performance measures directly correlate with the number of patterns identified. Groulx et al. (2017) indicates that Parallel coordinates approach makes identification of thresholds and correlation windows easy. Which can be valuable input for other regression efforts. Furthermore, the predictive proxy approach has a wide variety of applications to completion engineering and management. (Pankaj et al. 2018) Drilling and lifting in addition to completion also benefit from faster decision making according to Pankaj et al. (2018). Khodabakhsh, Ari, and Bakir (2017) agrees that the data analytic approach is applicable to the oil drilling industry in addition to oilwell completion practices. Subrahmanya et al. (2014) demonstrates that machine learning methods show considerable promise. Guevara et al. (2017) was able to outperform conventional techniques, such as kriging (a gaussian interpolation method). Future work in other domains was suggested by a variety of authors: vertical well logs (Guevara et al. 2017), geophysical (Luo 2018), fluid dynamics (Ezzatabadipour et al. 2017), exploration (Nikhalat-Jahromi and Jorge 2017), Remote sensing (Nikhalat-Jahromi and Jorge 2017), geobotany and geochemistry (Nikhalat-Jahromi and Jorge 2017). There is genuine excitement surrounding the application of machine learning and data analytics to various topics related to oil and gas production.

Feature Importance

Each author provides some insight into what features are most important for predicting oil production accurately. As a default position; High proppant (sand or other particulates used to prop open fractures in rock under high pressure), high-fluid (usually water, but sometimes gelling or cross-linking additives are included) completion designs described by Fulks et al. (2016) have shown success every basin. Fulks et al. (2016) goes on to say that degradable diversion (poly lactic acid diverts pressure but eventually degrades by bacteria) improves cluster efficiency and optimal lateral landing zone are important. In contrast, Temizel et al. (2015) lists in order of importance fracture half-length (which could be related to diversion), proppant amount, zone coverage (nearly equivalent to landing zone), and slurry volume (slurry is the combination of fluid with proppant). Ezzatabadipour et al. (2017) is more interested in two-phase flow patterns as a function of pipe condition. Luo (2018) purports that chemical ingredients are conducive to production and some, in fact, negatively impact gas production. (Luo 2018) Lopes and Jorge (2017) includes rock and fluid properties. While Guevara et al. (2017) expects to add well completion parameters to their model. Nikhalat-Jahromi and Jorge (2017) provides structural geology and reservoir properties (fault lines, water zones, etc.) as important features. Wang et al. (2016) focuses on trap and peel heights as well as gas flux, plans to add droplet size and gas-oil ratio to that list.

Simulation

Another common thread among recent literature is the idea of simulating new data from predictive models. Temizel et al. (2015) sees simulations as a tool to determine not only feature importance and significance, but also effect direction. Temizel et al. (2015) goes on to say that simulations can replace commercial fracture simulators. The approach described in Pankaj et al. (2018) creates a “parametric explosion of parameters”; enabling optimal well completion design to be determined much faster than with traditional methods, a matter of minutes instead of months.(Pankaj et al. 2018).Bozoev and Demidova (2016) agrees that simulation tends to be an appropriate approach to choosing the optimum completion of the wellbore. Like-wise, Ezzatabadipour et al. (2017) notes that investigations could be improved by exhaustive searches. Wang et al. (2016) also observes that resorting to simulated data from very high-resolution numerical simulations is a good compromise. Liu et al. (2018) uses simulation to study early hydration stages of a cement slurry which effect the transmission of hydrostatic pressure.

Model Limitations

Though there is much fanfare surrounding data analytic approaches applied to oilwell completions, many authors note that there are limitations to such endeavors. Fulks et al. (2016) says that establishing baseline performance is necessary. Temizel et al. (2015) takes the stance that conclusions from data driven models are specific to the model used. Careful consideration is necessary for normalizing both performance measures and inputs. Before discernible conclusions, patterns detected must be reviewed to refine insights. (Groulx et al. 2017) Pankaj et al. (2018) notes on that a calibrated model is fundamental to a reliable prediction. Bozoev and Demidova (2016) has concern that analytical approaches do not take into account interference effects. Wang et al. (2016) adds that next stage numerical models must be able to handle mixtures of oil and gas in order to simulate realistically.

Conclusion from Literature

For centuries, oilwells have been drilled vertically with little care about completion. Today, horizontal drilling has changed that. literature shows completion parameters have a significant effect on oil production. And

that these parameters can be used to make accurate predictions. While there is some disagreement about the limitations and applicability of such models, most agree that a general data analysis workflow is applicable.

Regression Analysis

WellDatabase is an aggregator of publically available data related to oil and gas production. Many states have an entity responsible for collecting and enforcing reporting requirements which vary between jurisdictions. In Texas, for example, Texas Railroad Commission (TXRRC) regulates the oil and gas industry. Location and depth of well casing is tightly regulated. TXRCC requires all hydrocarbon production that leaves the well site to be reported monthly. In the state of Texas, oil wells are tested for once per year. This analysis will focus on two counties in west Texas, Midland County and Pecos County. Our goal is to model 12 month cumulative oil production as a function of publically available drilling and completion parameters summarized in the following table.

column	description	units
api	14 digit unique well identifier	
surfacelatitude	surface location	degrees
surfacelongitude	surface location	degrees
county	Texas county	either Midland or Pecos
producing_formation	geologic formation	e.g Wolfcamp, Spraberry
wellboreprofile	type of drilling	Vertical, Directional, or Horizontal
trueverticaldepth	depth of production casing	feet
laterallength	horizontal component of well	feet
totalbasewatervolume	volume of water used	gallons
totalproppantmass	mass of sand used	lbs
fluidsystem	type of fluid additives	slick-water, linear, or cross-link
surfactantpresent	wether or not surfactant was used	True or False
claycontrolpresent	weather or not clay control was used	True or False
acidtreatmentpresent	weather or not acid was used	True or False
choke_size_clean	open proportion of valve	0 is closed, 1 is open
oil	first 12 month cumulative oil production	bbl

Preprocessing

Each table of wellDatabase must be summarize and deduplicate by api number then joined to the header on api where wellboreprofile=HORIZONTAL. For this analysis, we are only interested in horizontal wells. It was noted that totalproppantmass has many zeros. This is not physically meaningful, so when totalproppant-concentration=0 we will set it to null instead. The python code for this processing is available on github. It is tempting to use the reported monthly oil production as is, but that would lead to a target variable that is highly auto-correlated; meaning that the previous month's production is a good estimate of production in the next month. For this reason, we will focus our effort on modeling cumulative oil production. At this point, let's consider the scatter plot matrix.

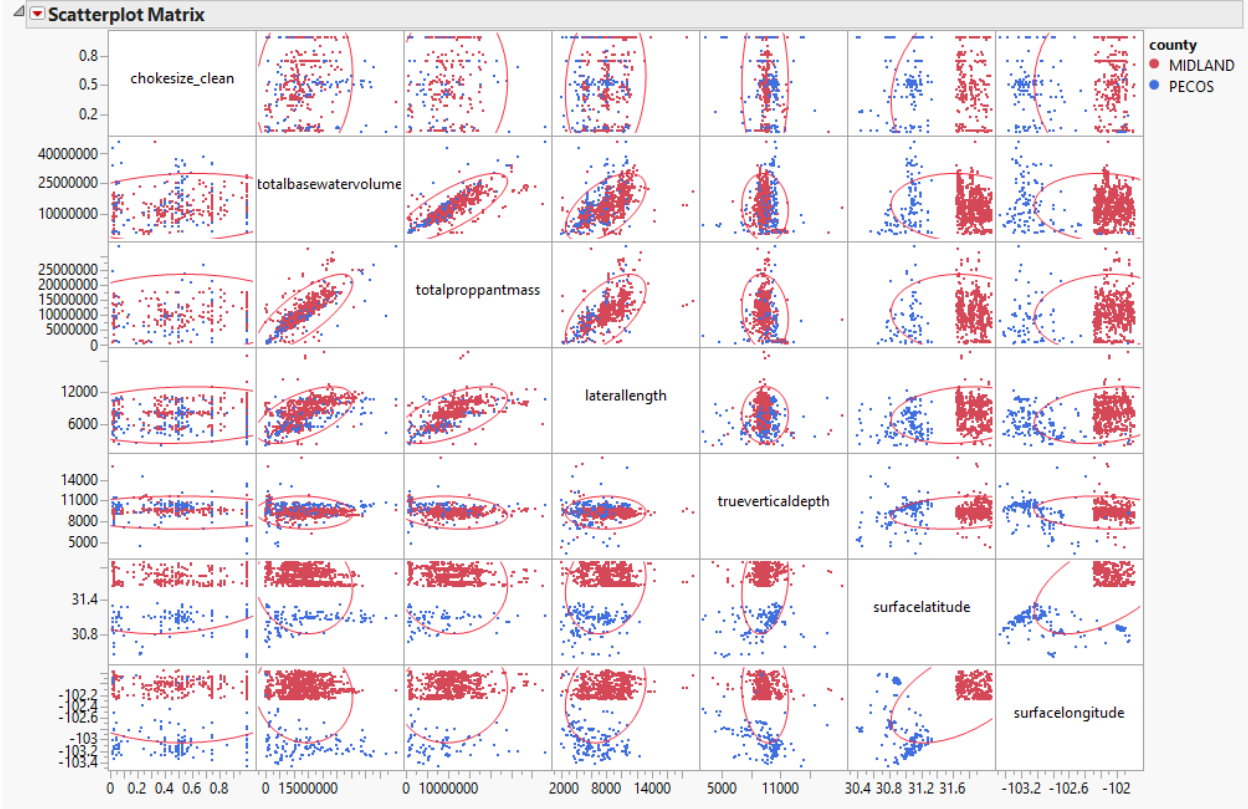
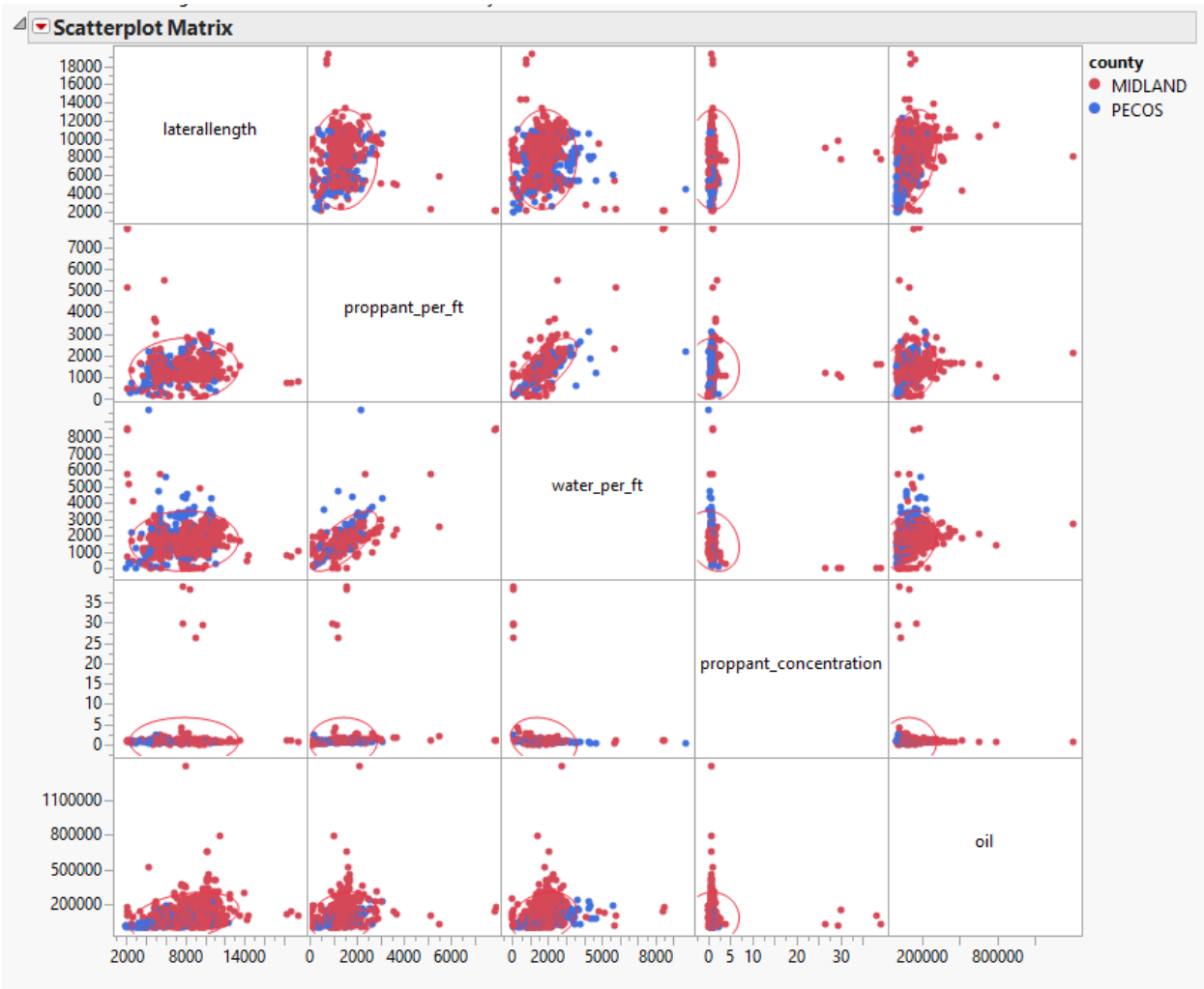


Figure 1: Predictors

From the scatterplot matrix, we can clearly see locational separation of the two counties in question. We can also see that totalproppantmass, totalbasewatervolume, and laterallength are highly correlated. To combat this, we will consider some new predictors which are scaled by length as well as the proppant concentration (proppant/water).

column	description	units
proppant_per_ft	totalproppantmass / laterlength	lbs/ft
water_per_ft	totalbasewatervolume / laterlength	gal/ft
proppant_concentration	totalproppantmass / totalbasewatervolume	lbs/gal

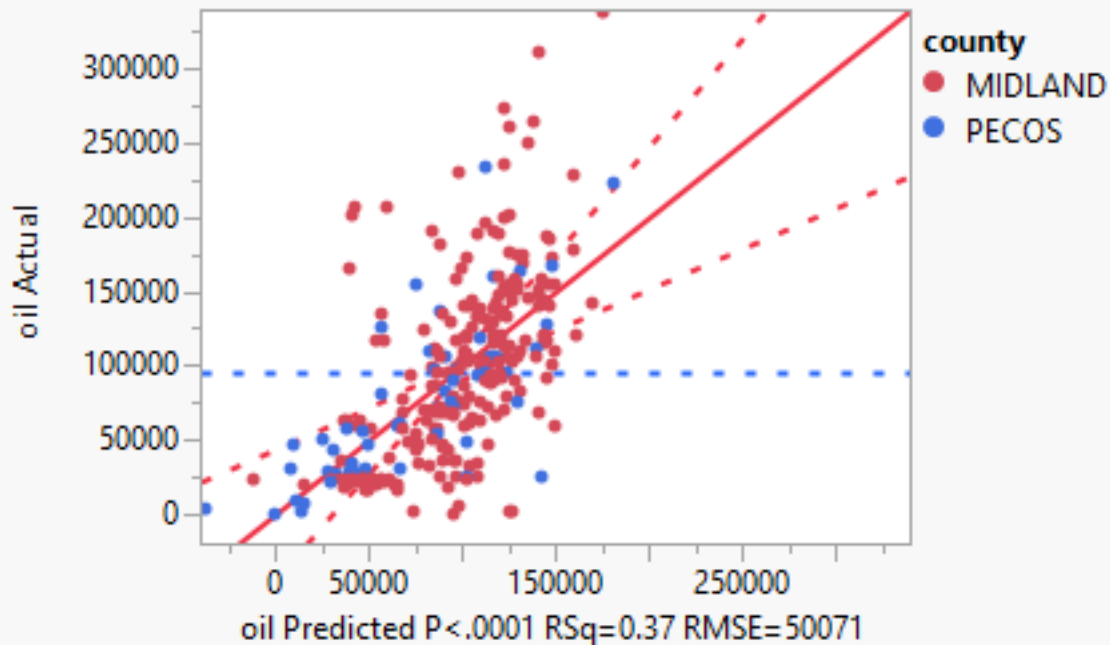


The newly calculated columns appear to be less correlated with each other and potentially better linear predictors of oil production.

Modelling

Even though it is expected that the relationship between the predictors and oil production is highly non-linear. Let's consider a naive main-effects-only model.

Actual by Predicted Plot



Summary of Fit

RSquare	0.371257
RSquare Adj	0.323443
Root Mean Square Error	50071.06
Mean of Response	95335.03
Observations (or Sum Wgts)	284

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	20	3.8934e+11	1.947e+10	7.7647
Error	263	6.5937e+11	2.5071e+9	Prob > F
C. Total	283	1.0487e+12		<.0001*

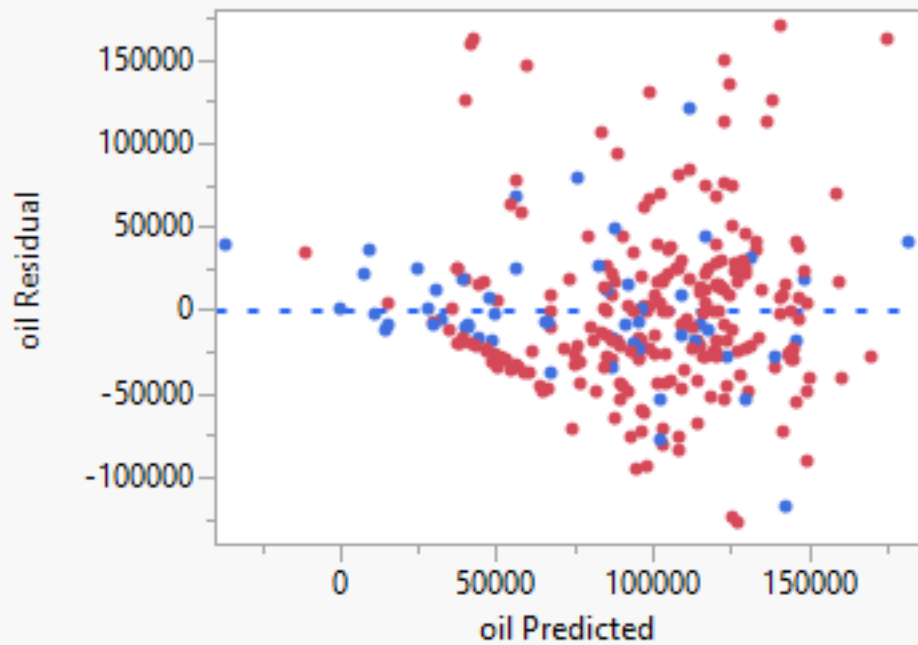
From the R^2 and R^2 adjusted values, this model explains about 30% of the variation in oil production. While not very accurate, this model could still be useful for industrial predictions.

Parameter Estimates				
Term	Estimate	Std Error	t Ratio	Prob> t
laterallength	8.6907153	1.578824	5.50	<.0001*
liftingmethod[gas lift]	-38792.62	8530.704	-4.55	<.0001*
proppant_per_ft	30.714412	12.38832	2.48	0.0138*
liftingmethod[flowing]	15346.755	6659.799	2.30	0.0220*
proppant_concentration	-17324.14	13292.28	-1.30	0.1936
trueverticaldepth	9.118752	7.012779	1.30	0.1946
surfacelatitude	38690.607	31030.9	1.25	0.2136
producing_formation[WOLFCAMP]	14177.336	11401.83	1.24	0.2148
surfactantpresent[False]	-13450	11888.4	-1.13	0.2589
surfacelongitude	28776.447	27558.32	1.04	0.2974
claycontrolpresent[False]	-4024.58	3897.555	-1.03	0.3027
liftingmethod[pumping]	-9200.774	9590.225	-0.96	0.3382
acidtreatmentpresent[False]	9415.5815	12237.91	0.77	0.4424
choke_size_clean	6455.9346	9341.646	0.69	0.4901
county[MIDLAND]	-13557.18	20674.68	-0.66	0.5126
producing_formation[SPRABERRY]	8268.3328	14090.7	0.59	0.5578
Intercept	1603815	2736993	0.59	0.5584
water_per_ft	5.2156039	9.87102	0.53	0.5977
producing_formation[DEAN]	7485.0438	18919.05	0.40	0.6927
fluidsystem[Linear]	1255.5922	4873.006	0.26	0.7969
fluidsystem[Crosslink]	946.19143	5145.189	0.18	0.8542

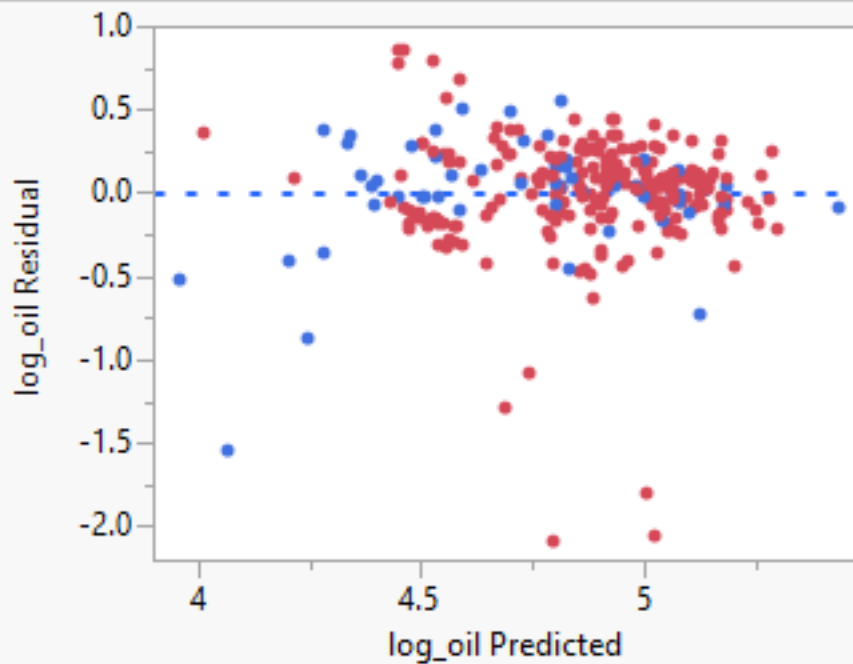
Based on an F statistic of < 0.0001 at least one coefficient is significantly different from zero, so we should proceed with attempting to improve this model. Based on t-statistics of each parameter we can see that relatively few coefficients are significantly different from zero. The most significant are laterallength, liftingmethod, and proppant_per_ft. A unit increase in proppant_per_ft is associated with an increase of 30 barrels of oil. Already we have performed about 20 hypothesis tests for significance.

Linear regression assumes homoscedasticity, meaning that the residuals are normally distributed with a mean of zero and constant variance i.e. $\epsilon \sim N(0, \sigma^2)$. From a plot of residual vs predicted oil, we can see a “megaphone” pattern which is evidence that this model violates the assumption of constant variance. In contrast, a similar main-effects only model predicting $\log(\text{oil})$ exhibits more constant variance in the residual.

Residual by Predicted Plot



Residual by Predicted Plot



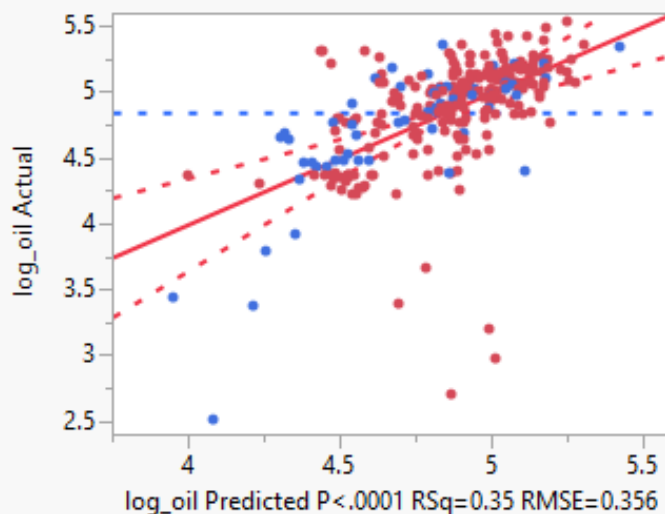
Going forward we will consider two models of $\log(\text{oil})$. The first is the main effects only model with the least significant terms removed, the second is a polynomial model up to degree two meaning that there will be a term to capture quadratic and interactive behavior. For the second model we will use step-wise regression to add terms until Bayesian information criteria (BIC) cannot be reduced by adding an additional term (forward selection).

Parameter Estimates				
Term	Estimate	Std Error	t Ratio	Prob> t
liftingmethod[gas lift]	-0.281893	0.061157	-4.61	<.0001*
laterallength	0.0000481	1.132e-5	4.25	<.0001*
liftingmethod[flowing]	0.1664615	0.047744	3.49	0.0006*
proppant_per_ft	0.0002384	8.881e-5	2.68	0.0077*
surfactantpresent[False]	-0.172537	0.085228	-2.02	0.0439*
surfacelatitude	0.4487527	0.22246	2.02	0.0447*
proppant_concentration	-0.162188	0.095292	-1.70	0.0899
acidtreatmentpresent[False]	0.1436435	0.087733	1.64	0.1028
claycontrolpresent[False]	-0.033396	0.027942	-1.20	0.2331
surfacelongitude	0.1957285	0.197565	0.99	0.3227
county[MIDLAND]	-0.132576	0.148217	-0.89	0.3719
liftingmethod[pumping]	-0.06037	0.068752	-0.88	0.3807
choke_size_clean	-0.053843	0.06697	-0.80	0.4221
trueverticaldepth	4.0325e-5	5.027e-5	0.80	0.4232
producing_formation[WOLFCAMP]	0.0574709	0.08174	0.70	0.4826
fluidsystem[Crosslink]	0.0249543	0.036886	0.68	0.4993
producing_formation[DEAN]	0.0891988	0.135631	0.66	0.5113
producing_formation[SPRABERRY]	-0.050675	0.101016	-0.50	0.6163
Intercept	9.6187998	19.62147	0.49	0.6244
fluidsystem[Linear]	0.0158116	0.034935	0.45	0.6512
water_per_ft	0.0000161	7.077e-5	0.23	0.8201

Removing water_per_ft, fluidsystem, and producing_formation results in the following reduced features main-effects only model.

Reduced Main Effects Model

Actual by Predicted Plot



Summary of Fit

RSquare	0.345803
RSquare Adj	0.312007
Root Mean Square Error	0.356002
Mean of Response	4.845055
Observations (or Sum Wgts)	286

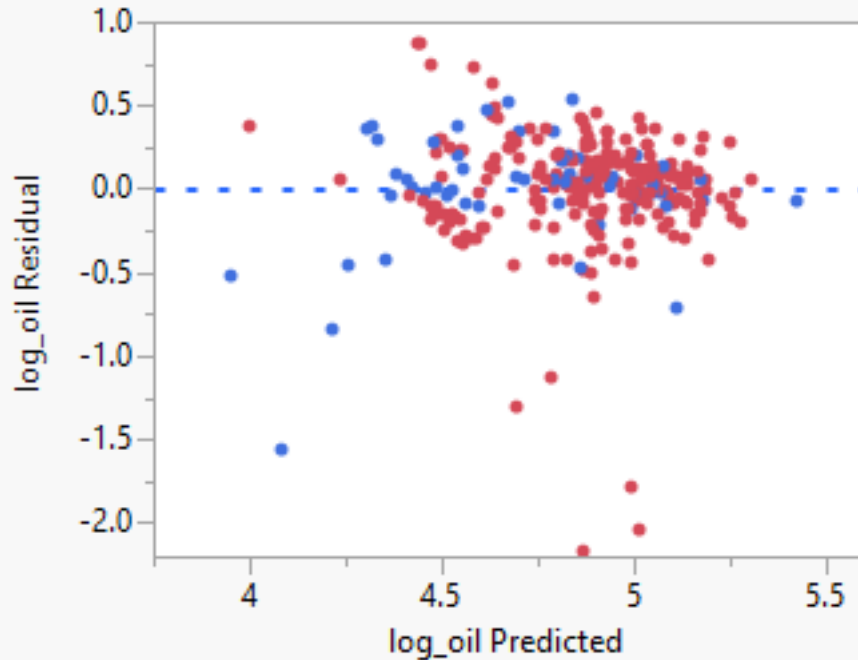
Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	14	18.154943	1.29678	10.2320
Error	271	34.345867	0.12674	Prob > F
C. Total	285	52.500810		<.0001*

Parameter Estimates

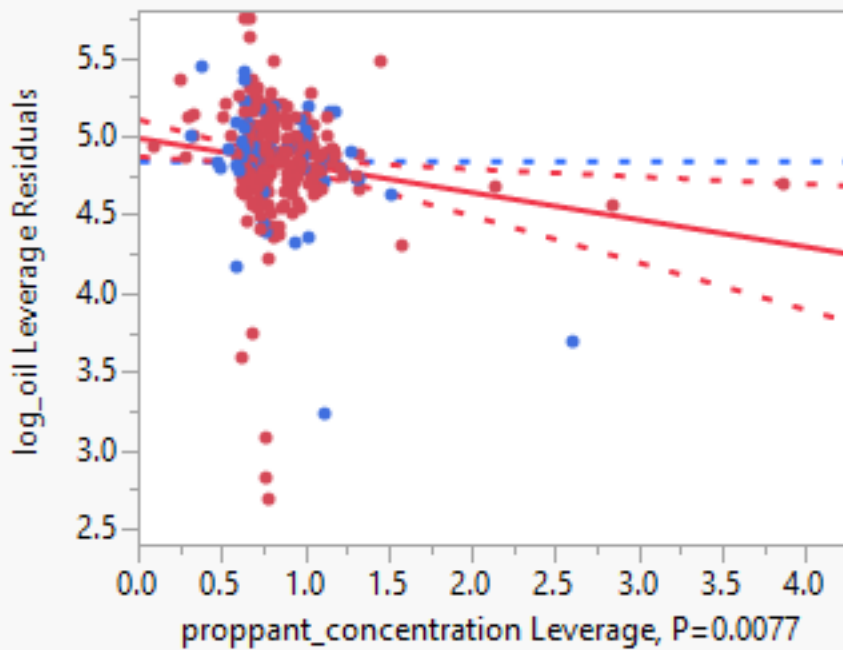
Term	Estimate	Std Error	t Ratio	Prob> t
proppant_per_ft	0.0002514	4.845e-5	5.19	<.0001*
laterallength	4.9032e-5	0.000011	4.47	<.0001*
liftingmethod[gas lift]	-0.263739	0.059248	-4.45	<.0001*
liftingmethod[flowing]	0.1586975	0.046738	3.40	0.0008*
proppant_concentration	-0.173076	0.064455	-2.69	0.0077*
surfactantpresent[False]	-0.145677	0.075947	-1.92	0.0561
trueverticaldepth	7.0316e-5	3.858e-5	1.82	0.0695
surfacelatitude	0.3211885	0.200757	1.60	0.1108
acidtreatmentpresent[False]	0.1027356	0.080674	1.27	0.2039
surfacelongitude	0.2326756	0.186441	1.25	0.2131
claycontrolpresent[False]	-0.031593	0.025915	-1.22	0.2239
liftingmethod[pumping]	-0.076637	0.066225	-1.16	0.2482
Intercept	17.195765	17.81376	0.97	0.3353
county[MIDLAND]	-0.11496	0.142252	-0.81	0.4197
choke_size_clean	-0.050029	0.065329	-0.77	0.4445

Residual by Predicted Plot

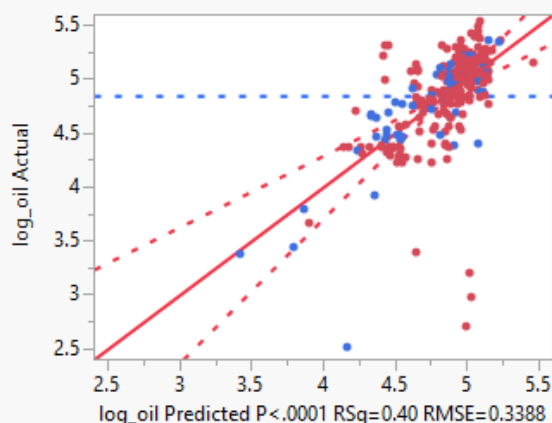


This model explains about 30% of the variance in $\log(\text{oil})$. Several main effects are still insignificant, but each coefficient is easily interpreted. For example, an increase of 1 lb/gal in proppant concentration is associated with a decrease 0.17 in $\log(\text{oil})$ or $\$ 10^{\{0.17\}} = 1.48\$$ barrels of oil

Leverage Plot



Actual by Predicted Plot



Summary of Fit

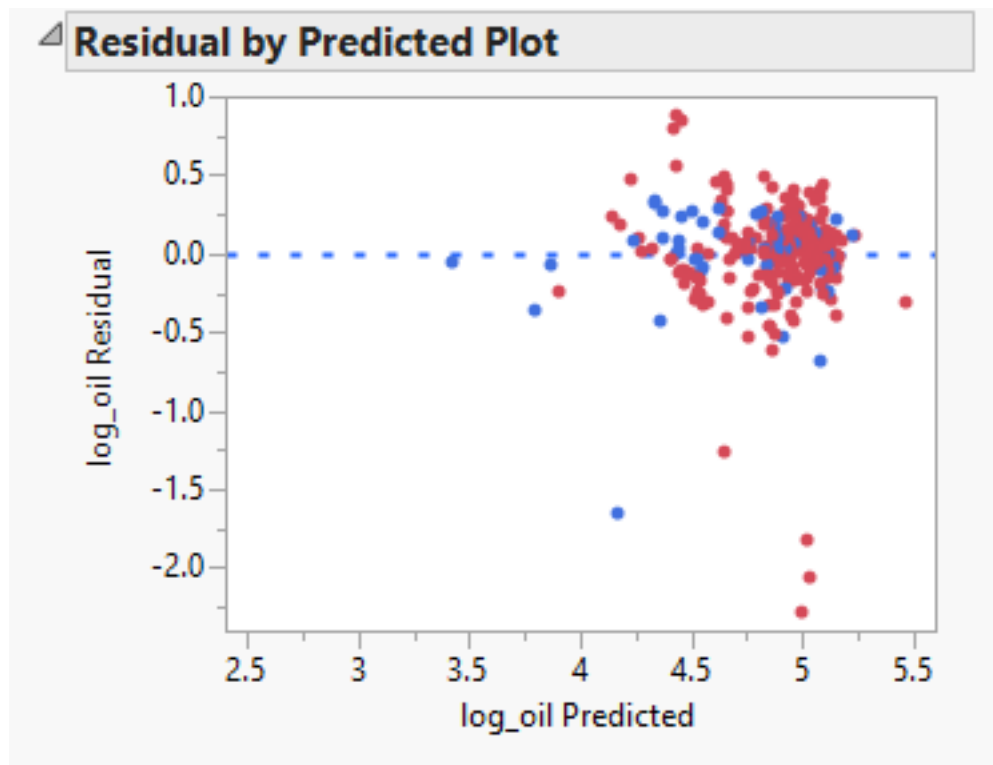
RSquare	0.403289
RSquare Adj	0.37706
Root Mean Square Error	0.338754
Mean of Response	4.845055
Observations (or Sum Wgts)	286

Analysis of Variance

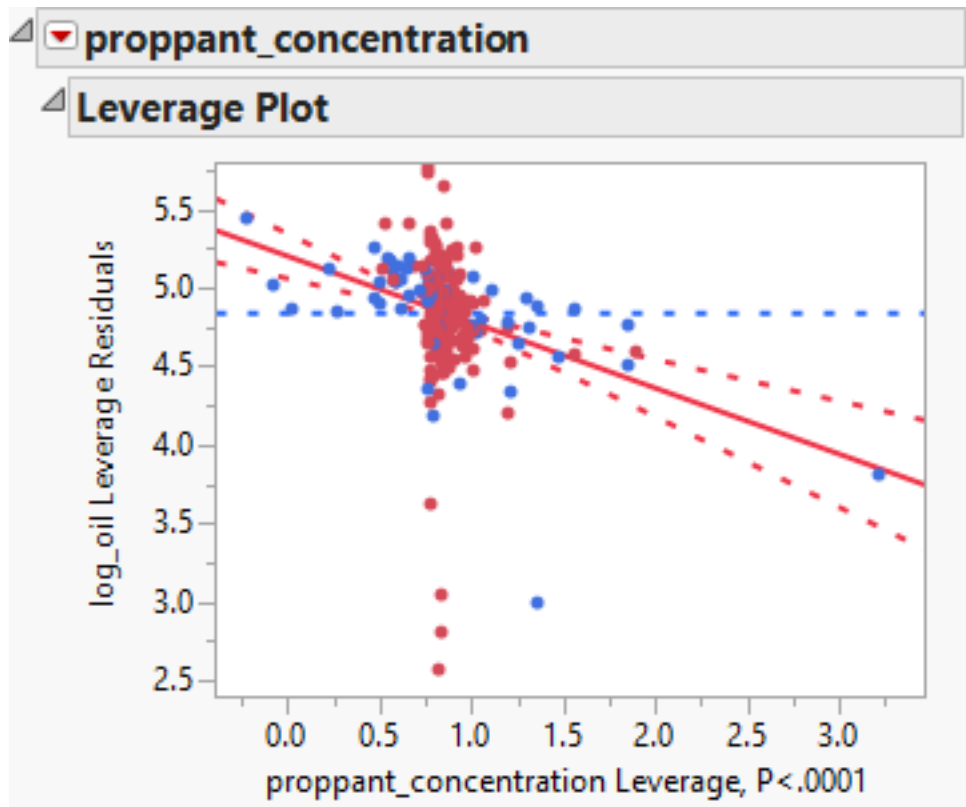
Source	DF	Sum of Squares	Mean Square	F Ratio
Model	12	21.172980	1.76442	15.3756
Error	273	31.327830	0.11475	Prob > F
C. Total	285	52.500810		<.0001*

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	3.7311663	0.198633	18.78	<.0001*
choke_size_clean	0.6978585	0.199452	3.50	0.0005*
liftingmethod{gas lift&pumping-steam injection&flowing}	-0.128707	0.027211	-4.73	<.0001*
liftingmethod{gas lift-pumping}	-0.107078	0.044469	-2.41	0.0167*
surfactantpresent[False]	-0.2851	0.078292	-3.64	0.0003*
laterallength	3.9412e-5	0.00001	3.95	<.0001*
county[MIDLAND]	0.1175034	0.027697	4.24	<.0001*
proppant_per_ft	0.000313	0.000064	4.90	<.0001*
proppant_concentration	-0.41953	0.079806	-5.26	<.0001*
(choke_size_clean-0.50187)*surfactantpresent[False]	0.7656701	0.19932	3.84	0.0002*
liftingmethod{gas lift-pumping}*(laterallength-7823.68)	-9.061e-6	2.441e-5	-0.37	0.7107
liftingmethod{gas lift&pumping-steam injection&flowing}*(proppant_per_ft-1312.99)	0.0001734	6.171e-5	2.81	0.0053*
county[MIDLAND]*(proppant_concentration-0.86184)	0.3135442	0.082393	3.81	0.0002*



This model explains nearly 40% of the variance, a marked improvement. Nearly all coefficient are significant, however some are much more difficult to interpret. In this model a change in proppant would effect more than 4 terms of our model in a non-linear manner.



Summary

It is possible to predict the first 12 months of oil production from publicly available drilling and completions information with a useful accuracy. In this analysis we performed a total of 54 hypothesis tests. However, refrained from setting any arbitrary significance level. Instead, we focused on the most and least significant terms as well as balancing predictive power with easy of interpretation. While complex a polynomial model explains more variation, a main-effects only model is equally useful and much easier to interpret. Proppant amount per foot and per gallon was the most significant feature regardless of model selection.

References

- Bozoev, A M, and E A Demidova. 2016. "Selection of the Optimal Completion of Horizontal Wells with Multi-Stage Hydraulic Fracturing of the Low-Permeable Formation, Field c." *IOP Conference Series: Earth and Environmental Science* 33 (1): 012035. <http://stacks.iop.org/1755-1315/33/i=1/a=012035>.
- Ezzatabadipour, Mohammadmehdi, Parth Singh, Melvin Deloyd Robinson, Pablo Guillen-Rondon, and Carlos Torres. 2017. "Deep Learning as a Tool to Predict Flow Patterns in Two-Phase Flow." *CoRR* abs/1705.07117. <http://arxiv.org/abs/1705.07117>.
- Fulks, Robert, Simon Hughes, Ingo Geldmacher, and others. 2016. "Optimizing Unconventional Completions-an Integrated Approach." In *Abu Dhabi International Petroleum Exhibition & Conference*. Society of Petroleum Engineers.
- Groulx, Bertrand, Jim Gouveia, Don Chenery, and others. 2017. "Multivariate Analysis Using Advanced Probabilistic Techniques for Completion Design Optimization." In *SPE Unconventional Resources Conference*. Society of Petroleum Engineers.
- Guevara, Jorge, Matthias Kormaksson, Bianca Zadrozny, Ligang Lu, John Tolle, Tyler Croft, Mingqi Wu, Jan Limbeck, and Detlef Hohl. 2017. "A Data-Driven Workflow for Predicting Horizontal Well Production Using Vertical Well Logs." *CoRR* abs/1705.06556. <http://arxiv.org/abs/1705.06556>.
- Holdaway, Keith R, Moray L Laing, and others. 2015. "Drilling and Completion Optimization in Unconventional Reservoirs with Data-Driven Models." In *SPE Asia Pacific Unconventional Resources Conference and Exhibition*. Society of Petroleum Engineers.
- Khodabakhsh, Athar, Ismail Ari, and Mustafa Bakir. 2017. "Cloud-Based Fault Detection and Classification for Oil & Gas Industry." *CoRR* abs/1705.04583. <http://arxiv.org/abs/1705.04583>.
- Liu, Kaiqiang¹, chengxw@swpu.edu.cn², j-656@163.com³ Cheng Xiaowei¹, Xingguo² Zhang, Zaoyuan² Li, Jia¹ Zhuang, and Xiaoyang² Guo. 2018. "Relationship Between the Microstructure/Pore Structure of Oil-Well Cement and Hydrostatic Pressure." *Transport in Porous Media* 124 (2): 463–78.
- Lopes, Rui L., and Alípio Jorge. 2017. "Mind the Gap: A Well Log Data Analysis." *CoRR* abs/1705.03669. <http://arxiv.org/abs/1705.03669>.
- Luo, & Zhang, H. 2018. "Mining Fracfocus and Production Data for Efficacy of Fracturing Fluid Formulations." In *Society of Petroleum Engineers*. Society of Petroleum Engineers.
- Nikhalat-Jahromi, Hamed, and Alípio M. Jorge. 2017. "An Overview of Data Mining Applications in Oil and Gas Exploration: Structural Geology and Reservoir Property-Issues." *CoRR* abs/1705.06345. <http://arxiv.org/abs/1705.06345>.
- Pankaj, Piyush, Steve Geetan, Richard MacDonald, Priyavrat Shukla, Abhishek Sharma, Samir Menasria, Tobias Judd, and others. 2018. "Need for Speed: Data Analytics Coupled to Reservoir Characterization Fast Tracks Well Completion Optimization." In *SPE Canada Unconventional Resources Conference*. Society of Petroleum Engineers.
- Subrahmanya, Niranjana, Peng Xu, Amr El-Bakry, Carmon Reynolds, and others. 2014. "Advanced Machine

Learning Methods for Production Data Pattern Recognition.” In *SPE Intelligent Energy Conference & Exhibition*. Society of Petroleum Engineers.

Temizel, C, S Purwar, A Abdullayev, K Urrutia, Aditya Tiwari, and others. 2015. “Efficient Use of Data Analytics in Optimization of Hydraulic Fracturing in Unconventional Reservoirs.” In *Abu Dhabi International Petroleum Exhibition and Conference*. Society of Petroleum Engineers.

Wang, Shitao, Mohamed Iskandarani, Ashwanth Srinivasan, W. Carlisle Thacker, Justin Winokur, and Omar M. Knio. 2016. “Propagation of Uncertainty and Sensitivity Analysis in an Integral Oil-Gas Plume Model.” *Journal of Geophysical Research: Oceans* 121 (5): 3488–3501.