# Oil & Gas Case Study Project
## Group 90

## I. Description of the data

In this oil & gas case study, we are predicting "CumOil12Month" (dependent/response variable) or in another word: the amount of oil produced per well within the first 12 months of a wells productive life.

**Data**

In the dataset, we have 6098 observations with 23 variables, including the response variable. Independent variables include operator/company names (categorical variable: 65 unique values total), completion dates, reservoirs (categorical variable: 9 unique values total) and other quantitative variables. All independent variables could be categorized by three, whether it gives us information about 1) completion, 2) geology, and 3) well spacing.

1) In the completion group, the variables include operator, completion date, amount of proppant and fluid, and percent of components of the reservoir, as well as the total cost of the horizontal well (unit in millions of dollars).
2) In the geology group, the variables include reservoir, depth, porosity, and pressure of the reservoir.
3) In the well spacing group, the variables include horizontal and vertical distances to the nearest offset well (unit in Feet).

**Initial explorations**

As an exploratory analysis, we've plotted scatter plots and histograms to investigate distributions of each variable, created a correlogram to identify any correlation across variables, and time series plot for investigating trend (based on completion dates).

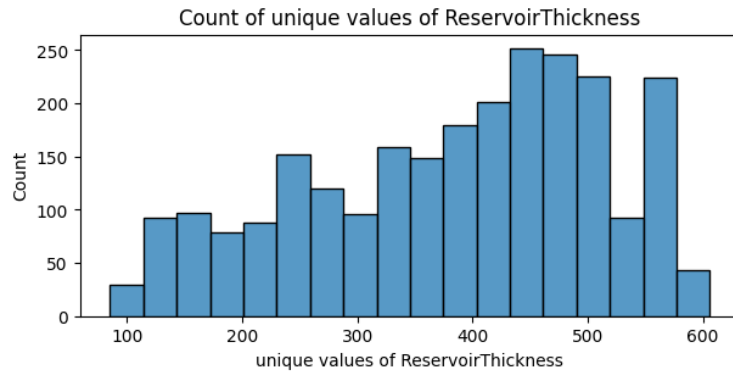We are only showing the most interesting visualizations here as we generate more plots and prints to understand the dataset.
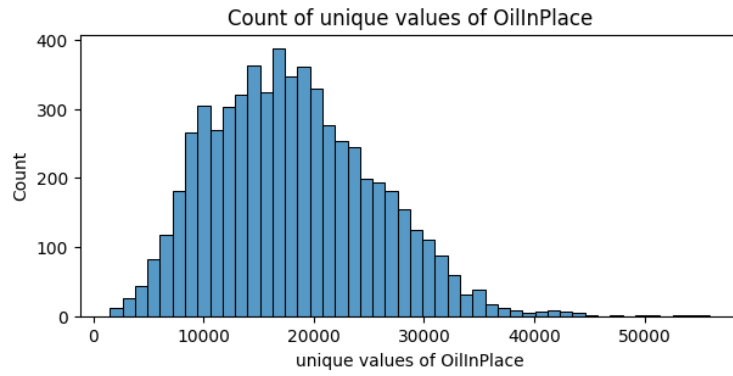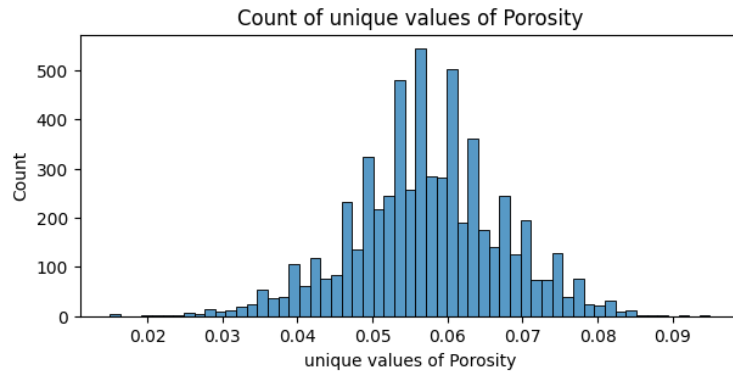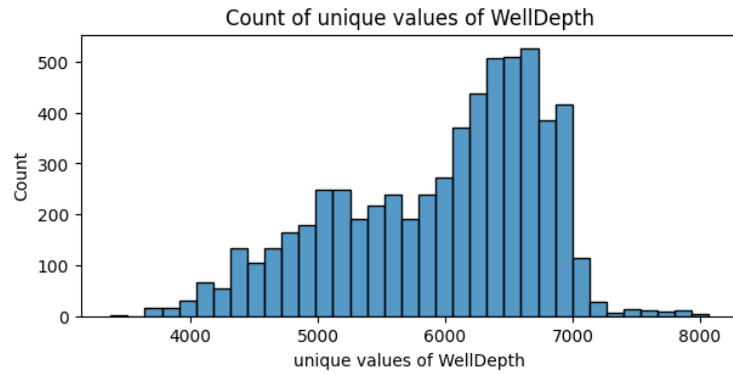
## II. Visualizations and captions

### a) Correlogram (%)

| | Operator | Reservoir | LateralLength_FT | ProppantIntensity_LBSPerFT | FluidIntensity_BBLPerFT | HzDistanceToNearestOffsetAtDrill | HzDistanceToNearestOffsetCurrent | VtDistanceToNearestOffsetCurrent | VtDistanceToNearestOffsetAtDrill | WellDepth | ReservoirThickness | OilInPlace | Porosity | ReservoirPressure | WaterSaturation | StructureDerivative | TotalOrganicCarbon | ClayVolume | CarbonateVolume | Maturity | TotalWellCost_USDMM | CumOil12Month |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Operator | -100 | 13 | 18 | -7 | -0 | -1 | -3 | 0 | -3 | -7 | 16 | 4 | 5 | 3 | 2 | 9 | -3 | -6 | 12 | -1 | 23 | 2 |
| Reservoir | 13 | 100 | -2 | -3 | -4 | 1 | 2 | 6 | -0 | 15 | | -27 | 13 | 59 | 27 | 7 | 2 | -32 | 51 | 8 | 5 | -16 |
| LateralLength_FT | 18 | -2 | 100 | 15 | 19 | -16 | -14 | -7 | -13 | -6 | 4 | 2 | 1 | 2 | -2 | 8 | -0 | -0 | 4 | 5 | 76 | 46 |
| ProppantIntensity_LBSPerFT | -7 | -3 | 15 | 100 | 54 | -9 | -6 | -8 | -10 | -1 | 7 | 3 | -2 | 0 | -3 | 3 | 2 | -1 | 1 | 8 | 57 | 28 |
| FluidIntensity_BBLPerFT | -0 | -4 | 19 | 54 | 100 | -6 | -4 | -5 | -8 | 13 | -10 | -1 | -10 | 9 | -4 | -4 | 0 | -9 | 4 | 2 | 45 | 33 |
| HzDistanceToNearestOffsetAtDrill | -1 | 1 | -16 | -9 | -6 | 100 | 78 | 23 | 51 | -3 | -8 | -6 | 1 | -5 | 4 | 6 | -0 | 8 | -6 | -10 | -12 | -7 |
| HzDistanceToNearestOffsetCurrent | -3 | 2 | -14 | -6 | -4 | 78 | 100 | 32 | 36 | -3 | -2 | -5 | 2 | -4 | 6 | 7 | 1 | 12 | -6 | -9 | -8 | -9 |
| VtDistanceToNearestOffsetCurrent | 0 | 6 | -7 | -8 | -5 | 23 | 32 | 100 | 61 | -6 | 8 | -3 | 9 | -2 | 11 | 4 | 2 | 9 | -2 | -6 | -7 | -11 |
| VtDistanceToNearestOffsetAtDrill | -3 | -0 | -13 | -10 | -8 | 51 | 36 | 61 | 100 | -7 | 5 | -0 | 4 | -7 | 3 | 5 | 1 | 10 | -4 | -9 | -15 | -9 |
| WellDepth | -7 | 15 | -6 | -1 | 13 | -3 | -3 | -6 | -7 | 100 | -67 | -12 | -21 | 73 | -7 | -40 | -3 | -32 | 10 | 34 | 7 | 25 |
| ReservoirThickness | 16 | | 4 | 7 | -10 | -8 | -2 | 8 | 5 | -67 | 100 | 60 | 38 | -49 | 25 | 29 | 5 | 29 | 21 | -9 | 8 | -21 |
| OilInPlace | 4 | -27 | 2 | 3 | -1 | -6 | -5 | -3 | -0 | -12 | 60 | 100 | 19 | -17 | -66 | -1 | -3 | 34 | -4 | 5 | 4 | 1 |
| Porosity | 5 | 13 | 1 | -2 | -10 | 1 | 2 | 9 | 4 | -21 | 38 | 19 | 100 | -6 | 34 | 10 | 3 | 50 | -42 | -4 | 4 | -11 |
| ReservoirPressure | 3 | 59 | 2 | 0 | 9 | -5 | -4 | -2 | -7 | 73 | -49 | -17 | -6 | 100 | 12 | -24 | 1 | -39 | 39 | 26 | 13 | 11 |
| WaterSaturation | 2 | 27 | -2 | -3 | -4 | 4 | 6 | 11 | 3 | -7 | 25 | -66 | 34 | 12 | 100 | 8 | 4 | -4 | -11 | -2 | 1 | -11 |
| StructureDerivative | 9 | 7 | 8 | 3 | -4 | 6 | 7 | 4 | 5 | -40 | 29 | -1 | 10 | -24 | 8 | 100 | 5 | 12 | 2 | -24 | 5 | -11 |
| TotalOrganicCarbon | -3 | 2 | -0 | 2 | 0 | -0 | 1 | 2 | 1 | -3 | 5 | -3 | 3 | 1 | 4 | 5 | 100 | 2 | 1 | -3 | 2 | -3 |
| ClayVolume | -6 | -32 | -0 | -1 | -9 | 8 | 12 | 9 | 10 | -32 | 29 | 34 | 50 | -39 | -4 | 12 | 2 | 100 | -70 | -17 | -5 | -9 |
| CarbonateVolume | 12 | 51 | 4 | 1 | 4 | -6 | -6 | -2 | -4 | 10 | 21 | -4 | -42 | 39 | -11 | 2 | 1 | -70 | 100 | 15 | 9 | -5 |
| Maturity | -1 | 8 | 5 | 8 | 2 | -10 | -9 | -6 | -9 | 34 | -9 | 5 | -4 | 26 | -2 | -24 | -3 | -17 | 15 | 100 | 9 | 14 |
| TotalWellCost_USDMM | 23 | 5 | 76 | 57 | 45 | -12 | -8 | -7 | -15 | 7 | 8 | 4 | 4 | 13 | 1 | 5 | 2 | -5 | 9 | 9 | 100 | 46 |
| CumOil12Month | 2 | -16 | 46 | 28 | 33 | -7 | -9 | -11 | -9 | 25 | -21 | 1 | -11 | 11 | -11 | -11 | -3 | -9 | -5 | 14 | 46 | 100 |

There are some variables showing high correlation (e.g., reservoir thickness & well depth, clayvolume & carbonate volume), so we would very much need models with the capacity to deal with overfitting.
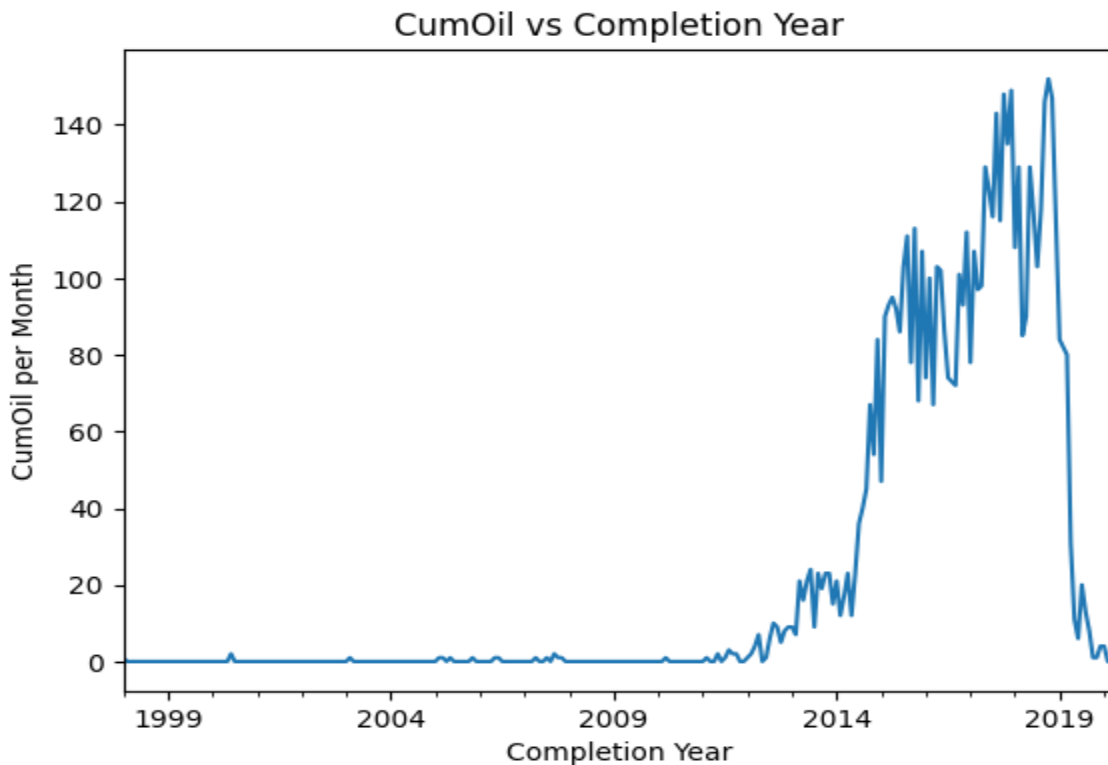
**b) Histograms of some of the key variables**


Count of unique values of WellDepth


Count of unique values of Porosity


Count of unique values of OilInPlace


Count of unique values of ReservoirThickness

Count of unique values of HzDistanceToNearestOffsetCurrent

We observed a wide range of distributions across variables. Some look normally distributed, and some are not. We may need to apply transformation for some variables. We even have variables with concentrated density in the minimum and the maximum of the value range.

**c) Completion date trend (# of obs per month)**


CumOil vs Completion Year

Most of the observations are concentrated in the period of 2014-2019. We observe very sporadic occurrences before 2011. We might have to truncate data, or be strategic when choosing a train/test set.

## III. Revised project question

We would like to develop a model to best predict the amount of oil produced per well within the first 12 months of a wells productive life ("CumOil12Month") from the given predictors using the given data. In this case the project question stays the same after the exploration.

## IIII. Baseline model

We train a multilinear regression model as a baseline model. In the first place, missing values in the dataset was an issue. We had to do a mode (which is not good for qualitative variables and we'll fix that later) imputation on both the train and the test dataset in order to perform the regression. For each column we impute missing values with the mode of observations.

Below is the output of our baseline model:

| Dep. Variable: | CumOil12Month | R-squared (uncentered): | 0.853 |
|---|---|---|---|
| Model: | OLS | Adj. R-squared (uncentered): | 0.852 |
| Method: | Least Squares | F-statistic: | 1408. |
| Date: | Sat, 03 Dec 2022 | Prob (F-statistic): | 0.00 |
| Time: | 03:38:35 | Log-Likelihood: | -60074. |
| No. Observations: | 4878 | AIC: | 1.202e+05 |
| Df Residuals: | 4858 | BIC: | 1.203e+05 |
| Df Model: | 20 | | |
| Covariance Type: | nonrobust | | |

We can see that the adjusted R-squared of the model is 85.2%. The adjusted R squared is a measure of goodness of fit of the model, in-sample.

We also try to look at the correlation between the explanatory variables. The correlogram shows the presence of explanatory variables with a high level of correlation which means we should expect some overfitting. Finally the baseline model's R-squared on the test set is around 38.20%, as opposed to 85.3% on the train set.

**IV. References**

1. Alkhammash EH. An Optimized Gradient Boosting Model by Genetic Algorithm for Forecasting Crude Oil Production. *Energies*. 2022; 15(17):6416.
https://doi.org/10.3390/en15176416

2. Ibrahim NM, Alharbi AA, Alzahrani TA, et al. Well Performance Classification and Prediction: Deep Learning and Machine Learning Long Term Regression Experiments on Oil, Gas, and Water Production. *Sensors (Basel)*. 2022;22(14):5326. Published 2022 Jul 16.
https://doi.org/10.3390/s22145326