

Oil and Gas Decisions

Christopher Miller, Dylan DeJean

March 3rd, 2023

Contents

Project Description	1
Introduction	1
Research Questions	2
Variables	2
Exploratory Data Analysis (EDA)	2
Data Cleaning and Outliers	2
Univariate Density Plots for Water and OilGas	4
Multivarite Plot	4
Statistical Analysis	5
Question 1 and 2 (Probit Models):	5
Logit Model	6
Question 3 (Failure-Time):	6
Technical Appendix A (R Script)	7
Technical Appendix B (Assumptions and Results)	12

Project Description

Introduction

In the midst of the COVID-19 pandemic, the oil and gas industry took a hit due to market demand plummeting. To keep profits as high as possible, oil and gas companies began to cut production for certain wells to either save resources for a period of higher prices or to save money by shutting in a company's least profitable wells. For this analysis, a well is identified as shut-in when its production during the current month is less than half of the production of the last month. A well is considered reopened when its production rises to at least 75% of the production from the month before it was shut-in. The data for this analysis consists of horizontal wells in the Bakken field of North Dakota owned by the companies XTO and Whiting. We also have data for horizontal wells in the Marcellus field in Pennsylvania owned by EQT. All data is from 2020.

Table 1
Variables

Variable Name	Type	Description
Company	Identification	Company Identification (EQT,WHIT,XTO)
WellName	Identification	Name of Specific Well within the Company
County	Explanatory	County that the Well is Located
Oil	Explanatory	Amount of Oil that is Produced in a Given Month
Wtr	Explanatory	Amount of Water that is Produced in a Given Month
Gas	Explanatory	Amount of Gas that is Produced in a Given Month
OilGas	Explanatory	Oil + Gas in Given Month
month	Identification	Month
shutin	Response	Well is Shut-In (1), Well is not Shut-In (2)

For Penn State’s Department of Energy and Mineral Engineering, we want to model shut-ins using both a probit model for the Bakken wells and a survival analysis model for all wells that we are interested in. For both models, the variables of interest are well size, waste water or other liquids produced, and geographical region of the wells. We also want to test if the models should be based on both XTO and Whiting, or if they should be separate models.

Research Questions

Question 1: How does size of production, size of waste water produced and location of the well effect if the well was shut-in after February 2020, in all three companies?

Question 2: In the Bakken Fields of North Dakota should the separate models for the companies XTO and Whiting (WHIT) be combined into one model?

Question 3: How does size of production, size of waste water produced and location of the well effect the probability the well survived through the end of year based on how many months it was shut-in from March through December?

Variables

We will first look at the variables that will be possibly used in the analysis of the well data (Table 1). It is important to note that the shutin variable is calculated using the definition of shut-in defined from Arash Dahi and Andrew Kleit, “We define a well as being shut-in if its production is less than half of its production in latest previous month. We define a well as having return to production if its production rises to 75 percent of its level in the pre-shut-in month.”

Exploratory Data Analysis (EDA)

Data Cleaning and Outliers

First, in Figure 1, we are going to first look at the percentage of wells shut-in each month. The graph involves all months, but is filtered on only wells that are shut-in in the months March through December. This is because any well that is shut-in the months prior to March is shut-in because of other reasons than the pandemic, and this report is trying to isolate the effect of variables on Covid-19 related shut-ins.

Other data cleaning procedures is that only horizontal wells were selected and the variable OilGas was created to add the production of Oil and Gas together to get an overall production metric. In terms of data

collection outliers, there is no obvious mis-entries or missing values so there is nothing to deal with in that regard, when looking at the univariate and multivariate graphs we will explore if there are any unreasonable points.

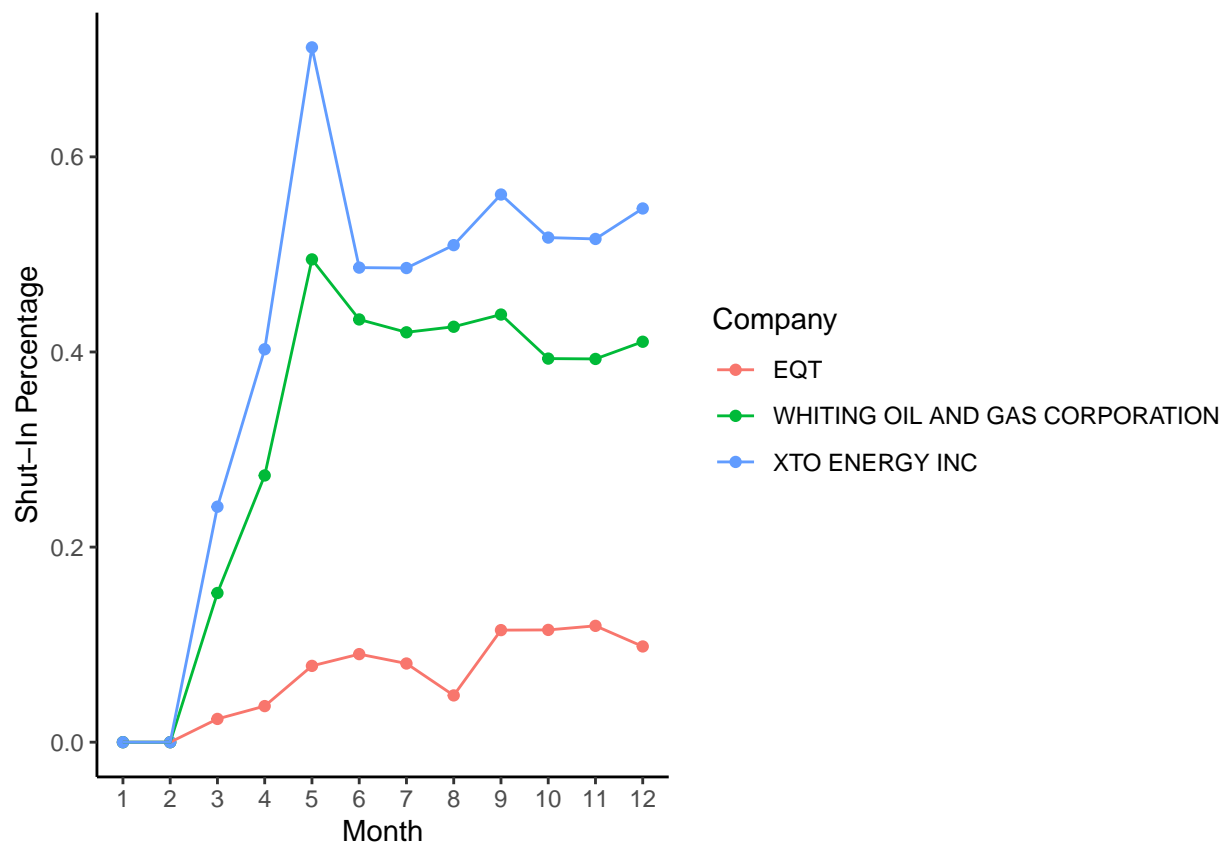


Figure 1: Shut-In Percentage Every Month

We can see that all of the WHIT and XTO, both of the companies in the North Dakota Bakken Fields, have an immediate and large response to the Covid-19 pandemic by shutting down many wells in March, April, and May. Both of those fields then begin to return some of the wells back to production, or possibly even open up more wells. There begins a small increase in shut-ins leading up into September, but then more wells begin to return to production as the year closes. As we can see for XTO, about 55% of wells are shut-in at the end of the year while WHIT has about 40% shut-in at the end of the year.

With the EQT company in Pennsylvania we see a different pattern emerge as the Covid pandemic effects increase throughout the year. As the months increase throughout the year 2020, EQT shut-ins wells gradually. Besides a dip going into September (opposite of the North Dakota Fields), we see a gradual increase into about 10% of shut-ins at the end of the year, an very large difference compared to the shut-ins of the North Dakota Fields.

Univariate Density Plots for Water and OilGas

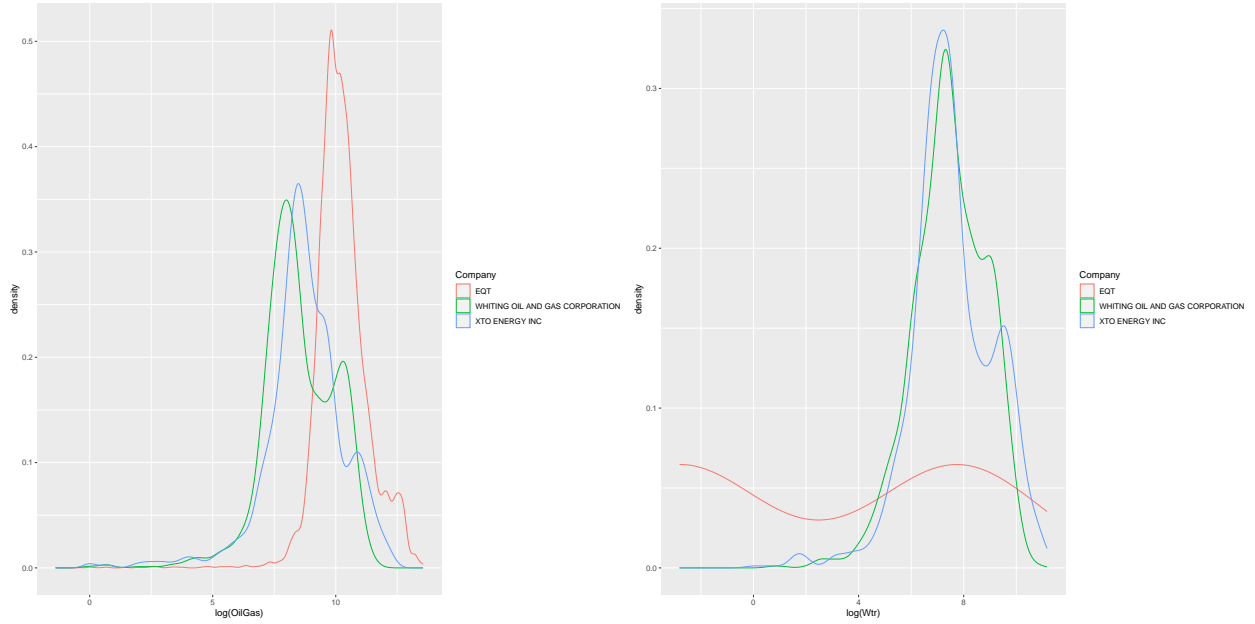


Figure 2: Water and OilGas Density Plots Separated by Companies

In Figure 2 the log of OilGas and Water is taken to lower the scale of the explanatory variables, and also get rid of all of the zero points where the well is shut-in. With OilGas we can see that WHIT and XTO are similar in terms of production with XTO having slightly higher production. EQT, the company in Pennsylvania, is much higher than WHIT and XTO, but it should be noted that EQT is only producing Gas, so that should be noted and considered in the analysis.

For Water, the North Dakota fields produce around the same water waste, but we can see there is a weird distribution for the Pennsylvania fields. This is because there are very few EQT wells with water waste creating this uneven odd distribution. The high density of zero water waste for Pennsylvania should be taken into account when modeling.

Multivariate Plot

We will now set the data into a model ready form where we will consider the observations to be a single well, NOT split into months. The response variable will be such that if the well was shut-in for any time-period past February we have (1 = shut-in) and if not (0 = not shut-in). We want measures for size of production and water waste quantity, so we will take the max of those variables over the months of productions for the wells to get an overall indicator of those variables for each well. We can then see how those variables are related with the response variable with boxplots.

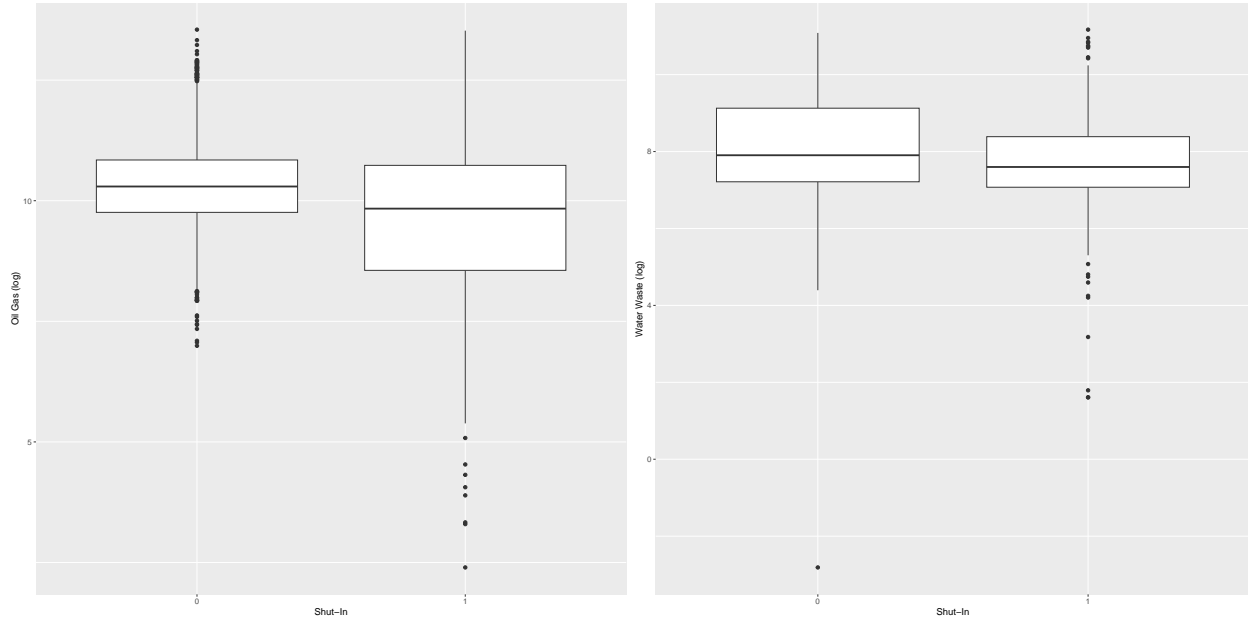


Figure 3: Shut-In vs. OilGas, Water Waste

Again we are working on the log scale because of the size of OilGas and Wtr. We can see that if a well was shut-in there is a larger variance of OilGas (size of production) and its median is slightly lower than well's that were never shut-in. In terms of water waste, the variance is similar, and the medians are similar if a well is shut-in or not.

We can see that there are some very large points and some very small points, but nothing that is justified enough to remove or transform, we could revisit this if model assumptions are not met.

It should be noted that when modeling we will not be working on the log scale for the explanatory variables because we model the companies separately and therefore scaling should be the same within each company.

Statistical Analysis

Question 1 and 2 (Probit Models):

First, we fit three models: one for Whiting horizontal wells only, one for XTO horizontal wells only, and one with both Whiting and XTO wells. Each model has three predictors: well size, the county where the well is located, and the amount of water and other liquid waste produced by the well. For well size, we decided that the best way to measure this would be to use the maximum oil and gas production for that year.

Unfortunately, due to the fact that the sample sizes of each model differ, there is no great way to directly compare the accuracy of each model to see if there should be separate models for Whiting and XTO. Despite this, we felt that the best way to figure out the optimal model would be to fit each, and then test to see if each model is significantly better than their respective null model (no predictors). If the combined model wasn't accurate, we would check the individual models and vice versa. To do this, we use a likelihood ratio test. We also want to test if county is an important predictor in each model since p-values were relatively high for county.

For the WHIT and XTO models, the likelihood ratio tests (LRT) to determine if county is significant have p-values greater than a pre-determined alpha level of .05, therefore we assume county is not a significant predictor in that model. For the combined model the p-value is less than .05, therefore county is significant

in that model. When testing each model compared to their respective null models, we find that the WHIT and combined models are highly significant. Since the XTO model is not an improvement on the null model, we will proceed with the combined model which explains both companies' well shut-in probability the best.

The final combined probit model is as follows:

$$P(Y = 1|WellSize, CountyDIV, CountyDUN, CountyMCK, CountyMTL, CountyWIL, Wtr) = \Phi(1.458 - 0.00002WellSize - 0.018CountyDIV - 1.338CountyDUN - 1.307CountyMCK - 1.399CountyMTL - 1.310CountyWIL + 0.00003Wtr) \quad (1)$$

Note: The probability of Y=1 in this scenario is the probability of the well being shut-in. All county variables are indicator variables as in they are either 0 or 1 depending on which county the well is in.

Logit Model

We wanted to explore logit models with the same predictors to see if there was any improvement. We first perform the same likelihood ratio tests to determine whether or not to include county in each model, and then we want to compare the models against their null models to find which models are significant. We find that county is only significant in the combined model and not in the WHIT or XTO models. Comparing to their respective null models, we find again that the WHIT and the combined models are significant, but not the XTO model. For the same reasons as the probit model, we proceed with the combined model as the best of the three.

To compare the combined probit and logit models, we will be using AIC. AIC is a way of measuring prediction error of a model and a lower value implies a more accurate model. The combined probit and logit models have AIC values of 667.05 and 666.64 respectively. This means that there is little to no difference between the two models when it comes to prediction accuracy, meaning that either can be used to proceed. Given the context of the study, we feel that logistic regression fits better. It is more interpretable as the response is log(odds) and is specially used to measure a binary outcome (shut-in or not) compared to the probit model which is used most frequently for a range of categories.

The combined logit model is as follows:

$$Log(odds) = 1.458 - 0.00002Size - 0.018CountyDIV - 1.338CountyDUN - 1.307CountyMCK - 1.399CountyMTL - 1.310CountyWIL + 0.00003Wtr \quad (2)$$

Question 3 (Failure-Time):

Technical Appendix A (R Script)

R Script

```
#Front Matter
knitr::opts_chunk$set(echo = FALSE, fig.pos = "H", out.extra = "")
#Setup
library(readxl)
library(dplyr)
library(ggplot2)
library(MASS)
library(lmtest)
library(survival)
library(survminer)
library(kableExtra)
library(gridExtra)
library(knitr)
library(effects)
rm(list = ls())
welldata <- {}
year = 1

#Dakota Data
while (year <= 12){

  if (year < 10){
    str1 = as.character(year)
    str2 = "2020_0"
    dataread <- paste(str2,str1, ".xlsx", sep = "")
    data <- read_xlsx(dataread)
    data$month <- year
    welldata <- rbind(welldata, data)
  }
  else{
    str1 = as.character(year)
    str2 = "2020_"
    dataread <- paste(str2,str1,".xlsx", sep = "")
    data <- read_xlsx(dataread)
    data$month <- year
    welldata <- rbind(welldata, data)
  }
  year = year+1
}

dakotawelldata <- welldata %>%
  filter(Pool == "BAKKEN"& Company == "XTO ENERGY INC" | Pool == "BAKKEN"& Company == "WHITING OIL AND GAS")
dakotawelldata$Oil <- as.numeric(dakotawelldata$Oil)
dakotawelldata$Gas <- as.numeric(dakotawelldata$Gas)
dakotaHorizontal <- read_xlsx('horizontal.xlsx')
dakotaHorizontal <- dakotaHorizontal %>%
```



```

#shut-in % by month
shutinbymonth <- dakotawelldata %>%
  group_by(Company,month)%>%
  summarise(shutinpct = sum(shutin)/n())

ggplot(shutinbymonth,aes(x=month,y=shutinpct,color=Company))+
  geom_point()+
  geom_line() +
  scale_x_continuous(breaks = round(seq(min(shutinbymonth$month), max(shutinbymonth$month), by = 1),1))
  labs(x="Month",y="Shut-In Percentage")+
  theme_classic()
plot1 <- ggplot(dakotawelldata, aes(x=log(OilGas), color=Company)) +
  geom_density()
plot2 <- ggplot(dakotawelldata, aes(x=log(Wtr), color=Company)) +
  geom_density()
grid.arrange(plot1, plot2, ncol=2)
dakotawelldata.s <- dakotawelldata %>%
  group_by(Company,WellName,County) %>%
  summarise(size = max(OilGas),Wtr = max(Wtr),shutin = sum(shutin) )
dakotawelldata.s$shutin <- ifelse(dakotawelldata.s$shutin > 0, 1,0)

plot1 <- ggplot(dakotawelldata.s, aes(x = as.factor(shutin), y=log(size))) +
  geom_boxplot()+
  labs(x = "Shut-In",y="Oil Gas (log)")
plot2 <- ggplot(dakotawelldata.s, aes(x = as.factor(shutin), y=log(Wtr))) +
  geom_boxplot()+
  labs(x = "Shut-In",y="Water Waste (log)")
grid.arrange(plot1, plot2, ncol=2)
dakotawelldata.s$County <- as.factor(dakotawelldata.s$County)

whiting_model <- glm(shutin ~ size + County + Wtr,
  data = dakotawelldata[dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION", ],
  family = binomial(link = "probit"))
whiting_null <- glm(shutin ~ 1, data = dakotawelldata[dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION", ],
  family = binomial(link="probit"))
whiting_no_county <- glm(shutin ~ size + Wtr,
  data = dakotawelldata[dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION", ],
  family=binomial(link="probit"))
xto_model <- glm(shutin ~ size + County + Wtr,
  data = dakotawelldata[dakotawelldata.s$Company == "XTO ENERGY INC", ],
  family = binomial(link = "probit"))
xto_null <- glm(shutin ~ 1, data = dakotawelldata[dakotawelldata.s$Company == "XTO ENERGY INC", ],
  family = binomial(link="probit"))
xto_no_county <- glm(shutin ~ size + Wtr,
  data = dakotawelldata[dakotawelldata.s$Company == "XTO ENERGY INC", ],
  family=binomial(link="probit"))
combined_model <- glm(shutin ~ size + County + Wtr,
  data = dakotawelldata[dakotawelldata.s$Company == "XTO ENERGY INC" |
    dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION", ],
  family = binomial(link = "probit"))
combined_no_county <- glm(shutin ~ size + Wtr,
  data = dakotawelldata[dakotawelldata.s$Company == "XTO ENERGY INC" |

```

```

                                dakotawelldata.s$Company == "WHITING OIL AND GAS CORP
                                family = binomial(link = "probit"))
combined_null <- glm(shutin ~ 1,
                    data = dakotawelldata.s[dakotawelldata.s$Company == "XTO ENERGY INC" |
                                dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION",
                                family = binomial(link = "probit"))
lrtest(whiting_model, whiting_no_county)
lrtest(xto_model, xto_no_county)
lrtest(combined_model, combined_no_county)
lrtest(combined_model, combined_null)
lrtest(whiting_model, whiting_null)
lrtest(xto_model, xto_null)
combined.final <- glm(shutin ~ size + County + Wtr,
                    data = dakotawelldata.s[dakotawelldata.s$Company == "XTO ENERGY INC" |
                                dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION",
                                family = binomial(link = "probit"))

whiting_model_logit <- glm(shutin ~ size + County + Wtr,
                        data = dakotawelldata.s[dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION",
                        family = binomial(link = "logit"))
whiting_null_logit <- glm(shutin ~ 1, data = dakotawelldata.s[dakotawelldata.s$Company
                                == "WHITING OIL AND GAS CORPORATION", ],
                        family = binomial(link="logit"))
whiting_no_county_logit <- glm(shutin ~ size + Wtr,
                        data = dakotawelldata.s[dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION",
                        family=binomial(link="logit"))
xto_model_logit <- glm(shutin ~ size + County + Wtr,
                    data = dakotawelldata.s[dakotawelldata.s$Company == "XTO ENERGY INC", ],
                    family = binomial(link = "logit"))
xto_null_logit <- glm(shutin ~ 1, data = dakotawelldata.s[dakotawelldata.s$Company == "XTO ENERGY INC",
                    family = binomial(link="logit"))
xto_no_county_logit <- glm(shutin ~ size + Wtr,
                    data = dakotawelldata.s[dakotawelldata.s$Company == "XTO ENERGY INC", ],
                    family=binomial(link="logit"))
combined_model_logit <- glm(shutin ~ size + County + Wtr,
                    data = dakotawelldata.s[dakotawelldata.s$Company == "XTO ENERGY INC" |
                                dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION",
                    family = binomial(link = "logit"))

combined_no_county_logit <- glm(shutin ~ size + Wtr,
                        data = dakotawelldata.s[dakotawelldata.s$Company == "XTO ENERGY INC" |
                                dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION",
                        family = binomial(link = "logit"))
combined_null_logit <- glm(shutin ~ 1,
                        data = dakotawelldata.s[dakotawelldata.s$Company == "XTO ENERGY INC" |
                                dakotawelldata.s$Company == "WHITING OIL AND GAS CORPORATION",
                        family = binomial(link = "logit"))

lrtest(combined_model_logit, combined_no_county_logit)
lrtest(whiting_model_logit, whiting_no_county_logit)
lrtest(xto_model_logit, xto_no_county_logit)
lrtest(whiting_no_county_logit, whiting_null_logit)
lrtest(xto_no_county_logit, xto_null_logit)

```

```

lrtest(combined_model_logit, combined_null_logit)
#time-to-event analysis

dakotawelldata <- dakotawelldata %>%
  group_by(WellName)%>%
  mutate(monthmax = max(month))
dakotawelldata$status <- ifelse(dakotawelldata$monthmax == dakotawelldata$month & dakotawelldata$shutin
whittime <- dakotawelldata %>%
  group_by(Company,WellName,County)%>%
  filter(Company == "WHITING OIL AND GAS CORPORATION")%>%
  summarise(failure = sum(shutin),status = sum(status), OilGas = max(OilGas),Wtr = max(Wtr))
xtotime <- dakotawelldata %>%
  group_by(Company,WellName,County)%>%
  filter(Company == "XTO ENERGY INC" )%>%
  summarise(failure = sum(shutin),status = sum(status), OilGas = max(OilGas),Wtr = max(Wtr))
combinedtime <- dakotawelldata %>%
  group_by(Company,WellName,County)%>%
  filter(Company == "XTO ENERGY INC" | Company == "WHITING OIL AND GAS CORPORATION")%>%
  summarise(failure = sum(shutin),status = sum(status), OilGas = max(OilGas),Wtr = max(Wtr))
eqttime <- dakotawelldata %>%
  group_by(Company,WellName,County)%>%
  filter(Company != "XTO ENERGY INC" & Company != "WHITING OIL AND GAS CORPORATION")%>%
  summarise(failure = sum(shutin),status = sum(status), OilGas = max(OilGas),Wtr = max(Wtr))

res.cox.whit <- glm(formula = status ~ failure+OilGas+County+Wtr,
  family = binomial(link = "cloglog"),
  data = whittime)
res.cox.xto <- glm(formula = status ~ failure+OilGas+County+Wtr,
  family = binomial(link = "cloglog"),
  data = xtotime)
res.cox.combined <- glm(formula = status ~ failure+OilGas+County+Wtr,
  family = binomial(link = "cloglog"),
  data = combinedtime)
res.cox.eqt <- glm(formula = status ~ failure+OilGas+County+Wtr,
  family = binomial(link = "cloglog"),
  data = eqttime)
plot(allEffects(res.cox.whit))

# Reprinted code chunks used previously for analysis

```

Technical Appendix B (Assumptions and Results)