

# Hawaii Water Data Exploration

*Gaby Garcia*

4/3/2019

## Research Question and Rationale

Which of the parameters have a relationship with dissolved oxygen concentrations? Of these relevant parameters, which have the most significant effect on dissolved oxygen concentrations over time? 2. Do dissolved oxygen concentrations vary spatially across Hawaii? Based on the 373 sample locations, I will code them based on geographic region on the big island (ex. west, east, north, south).

What are the water quality parameters that have a significant effect on Dissolved oxygen concentrations? (multiple linear regression)

## Dataset Information

### Explatory Data Analysis and Wrangling

```
setwd("~/Desktop/Environmental Data Analytics/Environmental_Data_Analytics/Final Project")
HawaiiWater<-read.csv('HawaiiWaterDataProcessed.csv')
```

## Load Necessary Packages

```
library(tidyverse)
library(tidyr)
library(ggplot2)
library(GGally)
library(dplyr)
library(plyr)
library(lubridate)
library(viridis)
library(RColorBrewer)
library(colormap)
library(gridExtra)
library(corrplot)
library(nlme)
library(lsmeans)
library(multcompView)
library(trend)
library(mapview)
library(leaflet)
library(sf)
library(car)
library(stats)
library(wesanderson)
```

```
library(scales)  
library(extrafont)
```

Omit NA's from Data (GLM 12 lesson says to do so)

```
HawaiiWaterClean<- na.omit(HawaiiWater)
```

## Convert Station Number from Factor to Number

```
HawaiiWaterClean$Station.No<-as.numeric(HawaiiWaterClean$Station.No)  
HawaiiWaterCleanGrouped<-HawaiiWaterClean%>%group_by(Location)
```

## Filtering Data Set to Only Include Observations from Oahu

```
HawaiiWaterCleanOahu <-  
  HawaiiWaterClean %>%  
    dplyr::filter(Location == "Ala Moana Park, Ewa" | Location == "Ala Moana Park, Center" | Location == "A  
" | Location == "Tongg's" | Location == "Kawaikui Beach Park" | Location == "Kanenelu Beach" | Location == "Kal  
" | Location == "Kaluahole Beach" | Location == "Outrigger" | Location == "Halona Cove" | Location == "Kokee Be  
" | Location == "Turtle Bay" | Location == "Kaunala Beach" | Location == "Pupukea at Shark's Cove" | Location =  
" | Location == "Kualoa Sugar Mill Beach" | Location == "Swanzys Beach Park" | Location == "Makaua Beach" | Locat
```

## Recoding Geographical Observations to Regions in Oahu

Convert station number from number to factor

```
HawaiiWaterCleanOahu$Station.No<-as.factor(HawaiiWaterCleanOahu$Station.No)
```

## Use Revalue function to recode Station Numbers into Geographical Regions

```
library(plyr)
HawaiiWaterCleanOahu$Region<-revalue(HawaiiWaterCleanOahu$Station.No, c('28'='East', '39'='South', '40'='North'))
```

Determine Number of Observations for Each Region in Oahu

```
summary(HawaiiWaterCleanOahu$Region)
```

```
##   East South North West
##   999 2634 1002 969
```

# Exploratory Data Analysis

## Structure of Water Data

```
str(HawaiiWaterCleanOahu)
```

```
## 'data.frame': 5604 obs. of 17 variables:
##   $ Sample.No          : Factor w/ 23020 levels "CF01030501","CF01030502",...: 4477 10741 ...
##   $ Sampler            : Factor w/ 21 levels "", "CF", "CF ", ...: 5 9 18 8 19 9 8 5 19 9 ...
##   $ Lab.No              : Factor w/ 15097 levels "", "1", "10", "100", ...: 14982 13031 12810 ...
##   $ Station.No          : Factor w/ 118 levels "28", "39", "40", ...: 2 2 2 2 2 2 2 2 2 ...
##   $ Location            : Factor w/ 240 levels "Airport", "Ala Moana DH1", ...: 7 7 7 7 7 7 ...
##   $ Date                : Factor w/ 1133 levels "1/10/00", "1/10/01", ...: 226 169 149 1055 ...
##   $ Time                : Factor w/ 501 levels "", "1:00:00 AM", ...: 277 427 437 431 52 422 ...
##   $ Enterococcus        : num 400 20 10 2.3 6.3 2.3 0.3 1 2.7 2.3 ...
##   $ CP                  : num 4.8 3 0.2 1 0.2 1 0.2 0.2 0.2 1 ...
##   $ Temperature         : num 25.1 25.6 26.9 26.7 26.4 ...
##   $ Salinity             : num 35.2 35.3 34.6 35.3 35.1 ...
##   $ DO                  : num 6.45 5.57 5.46 4.89 6.24 6.62 5.97 4.42 4.8 6.08 ...
##   $ PercentSaturationDissolvedOxygen: num 94 84.4 84.3 75.5 95.6 98.3 88.3 65 71 89.5 ...
##   $ pH                  : num 8.15 8.02 8.05 8.1 8.2 8.15 8.06 8.1 7.99 8.05 ...
##   $ Turbidity            : num 14.13 28.4 9.57 5.64 2.76 ...
##   $ Remarks              : Factor w/ 5481 levels "", "", "CLEAR/SUNNY", "SWIMMERS", ...: 2981 ...
##   $ Region               : Factor w/ 4 levels "East", "South", ...: 2 2 2 2 2 2 2 2 2 ...
## - attr(*, "na.action")= 'omit' Named int 2 46 73 80 81 82 83 85 86 87 ...
## ..- attr(*, "names")= chr "2" "46" "73" "80" ...
```

## Summary of Water Data

```
summary(HawaiiWaterCleanOahu)
```

```
##      Sample.No       Sampler       Lab.No       Station.No
## CF01190609: 1    JD :1011    H038-06: 1    92     : 227
## CF01250601: 1    SM : 923    H047-06: 1    48     : 226
## CF02020601: 1    SN : 922    H069-06: 1    68     : 226
## CF02150601: 1    GH : 851    H084-06: 1    74     : 224
## CF03020601: 1    JM : 767    H105-06: 1    71     : 223
## CF03090602: 1    DM : 619    H119-06: 1    82     : 223
## (Other) :5598 (Other): 511 (Other):5598 (Other):4255
##      Location        Date        Time
## San Souci       : 227 12/13/04: 40 8:30:00 AM: 153
## Kailua Beach    : 226 11/18/04: 38 7:40:00 AM: 139
## Kuhio Beach     : 226 11/8/04 : 38 7:50:00 AM: 120
## Hanauma Bay     : 224 11/16/04: 37 7:35:00 AM: 119
## Waialae-Kahala Beach: 223 11/22/04: 37 7:25:00 AM: 118
## Waimanalo Beach : 223 11/4/04 : 37 7:20:00 AM: 116
## (Other) :4255 (Other) :5377 (Other) :4839
##      Enterococcus     CP        Temperature      Salinity
## Min.   : 0.30  Min.   : 0.200  Min.   :20.29  Min.   : 8.42
## 1st Qu.: 2.30  1st Qu.: 0.500  1st Qu.:24.05  1st Qu.:34.54
## Median : 3.30  Median : 1.000  Median :25.14  Median :34.95
## Mean   : 28.51  Mean   : 3.055  Mean   :25.04  Mean   :34.31
```

```

## 3rd Qu.: 10.00 3rd Qu.: 2.000 3rd Qu.:26.07 3rd Qu.:35.18
## Max. :22000.00 Max. :290.000 Max. :28.91 Max. :37.24
##
##          DO      PercentSaturationDissolvedOxygen      pH
##  Min.   :3.060   Min.   : 5.82                  Min.   :7.340
##  1st Qu.:5.470   1st Qu.: 83.10                 1st Qu.:8.020
##  Median :5.880   Median : 88.80                 Median :8.110
##  Mean   :5.829   Mean   : 87.38                 Mean   :8.103
##  3rd Qu.:6.210   3rd Qu.: 92.70                 3rd Qu.:8.190
##  Max.   :9.340   Max.   :134.60                 Max.   :8.800
##
##          Turbidity                    Remarks      Region
##  Min.   : 0.000                  :3014     East : 999
##  1st Qu.: 1.800    SUNNY, LIGHT BREEZE, SWIMMERS: 18     South:2634
##  Median : 3.440    Tide rising, choppy, no rain : 15     North:1002
##  Mean   : 6.331    SUNNY, LIGHT BREEZE             : 14     West : 969
##  3rd Qu.: 7.425    SUNNY                         : 13
##  Max.   :315.000   SWIMMERS                      : 13
##                                     (Other)           :2517

```

## Dimensions of Data

```
dim(HawaiiWaterCleanOahu)
```

```
## [1] 5604 17
```

## View First 10 Rows of Data Frame

```
head(HawaiiWaterCleanOahu, 10)
```

```

##      Sample.No Sampler Lab.No Station.No          Location      Date
## 1 GH11160401     GH 0846-04      39 Ala Moana Park, Ewa 11/16/04
## 2 JM10260602     JM 02221-06      39 Ala Moana Park, Ewa 10/26/06
## 3 SM10200502     SM 02120-05      39 Ala Moana Park, Ewa 10/20/05
## 4 JD09140602     JD 01935-06      39 Ala Moana Park, Ewa 9/14/06
## 5 SN05050505     SN 00907-05      39 Ala Moana Park, Ewa 5/5/05
## 6 JM05180601     JM 01043-06      39 Ala Moana Park, Ewa 5/18/06
## 7 JD03100505     JD 00473-05      39 Ala Moana Park, Ewa 3/10/05
## 8 GH11220401     GH 0922-04      39 Ala Moana Park, Ewa 11/22/04
## 9 SN01200505     SN 00113-05      39 Ala Moana Park, Ewa 1/20/05
## 10 JM03020601    JM 00431-06      39 Ala Moana Park, Ewa 3/2/06
##
##            Time Enterococcus CP Temperature Salinity DO
## 1 6:15:00 AM      400.0 4.8    25.14  35.17 6.45
## 2 8:45:00 AM      20.0 3.0    25.62  35.27 5.57
## 3 8:55:00 AM      10.0 0.2    26.91  34.59 5.46
## 4 8:49:00 AM       2.3 1.0    26.66  35.32 4.89
## 5 10:00:00 AM      6.3 0.2    26.40  35.09 6.24
## 6 8:40:00 AM       2.3 1.0    24.60  34.84 6.62
## 7 10:03:00 AM      0.3 0.2    24.71  34.40 5.97
## 8 6:00:00 AM       1.0 0.2    25.51  35.26 4.42
## 9 8:27:00 AM       2.7 0.2    24.20  35.41 4.80
## 10 8:30:00 AM      2.3 1.0    24.17  34.82 6.08

```

```

##      PercentSaturationDissolvedOxygen    pH Turbidity
## 1                      94.0 8.15     14.13
## 2                      84.4 8.02     28.40
## 3                      84.3 8.05      9.57
## 4                      75.5 8.10      5.64
## 5                      95.6 8.20      2.76
## 6                      98.3 8.15      2.18
## 7                      88.3 8.06      6.10
## 8                      65.0 8.10      3.31
## 9                      71.0 7.99      3.41
## 10                     89.5 8.05      3.78
##
##                                         Remarks
## 1           S. CON. 53.2, NO WIND, LITTLE CLOUDS, SMALL WAVES, LITTLE MURKY
## 2
## 3
## 4
## 5           SP COND: 53.1, PEOPLE SURFING, OVERCAST, VERY LITTLE WAVE ACTION, CLEAR WATER
## 6
## 7           CALM, SWIMMERS, WINDY
## 8           S. COND 53.4, NO WIND, OVERCAST, CALM CLEAN WATER
## 9   SP COND 53.7, DEBRIS ON BEACH, SUNNY, CLEAR SKIES, LIGHT WIND, LOW TIDE, CLEAR WATER
## 10
##      Region
## 1   South
## 2   South
## 3   South
## 4   South
## 5   South
## 6   South
## 7   South
## 8   South
## 9   South
## 10  South

```

## View all Column Names

```
colnames(HawaiiWaterCleanOahu)
```

```

## [1] "Sample.No"                      "Sampler"
## [3] "Lab.No"                          "Station.No"
## [5] "Location"                        "Date"
## [7] "Time"                            "Enterococcus"
## [9] "CP"                               "Temperature"
## [11] "Salinity"                         "DO"
## [13] "PercentSaturationDissolvedOxygen" "pH"
## [15] "Turbidity"                        "Remarks"
## [17] "Region"

```

## Change Date from Factor to Date Object

```
HawaiiWaterCleanOahu$Date<-as.Date(HawaiiWaterCleanOahu$Date, format = "%m/%d/%y")
```

## Add a Week, Month, and Year Column to Dataframe Using Mutate Function

```
HawaiiWaterCleanOahu<-mutate(HawaiiWaterCleanOahu, Week = week(Date))
HawaiiWaterCleanOahu<- mutate(HawaiiWaterCleanOahu, Month = month(Date))
HawaiiWaterCleanOahu<- mutate(HawaiiWaterCleanOahu, Year = year(Date))
```

## Data Visualization

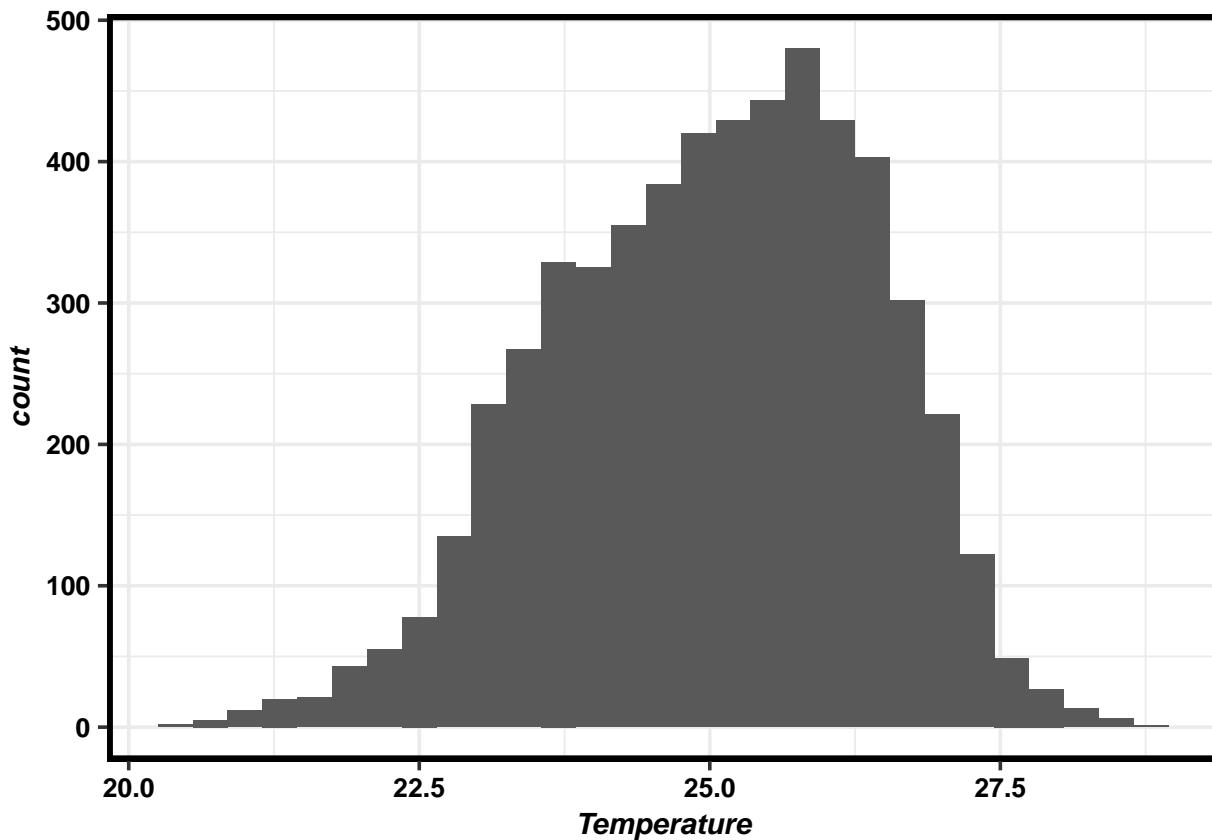
### Set GGPlot Theme

```
gabytheme <- theme_bw(base_size = 14) +
  theme(plot.title=element_text(face="bold", size="15", color="Indianred4", hjust=0.5),
         axis.title=element_text(face="bold.italic", size=11, color="black"),
         axis.text = element_text(face="bold", size=10, color = "black"),
         panel.background=element_rect(fill="white", color="darkblue"),
         panel.border = element_rect(color = "black", size = 2),
         legend.position = "top", legend.background = element_rect(fill="white", color="black"),
         legend.key = element_rect(fill="transparent", color="NA"))
theme_set(gabytheme)
```

**Examine distributions of continuous variables: Temperature, pH, Dissolved Oxygen Concentrations, Salinity, Turbidity, and Enterococcus Concentrations**

### Temperature

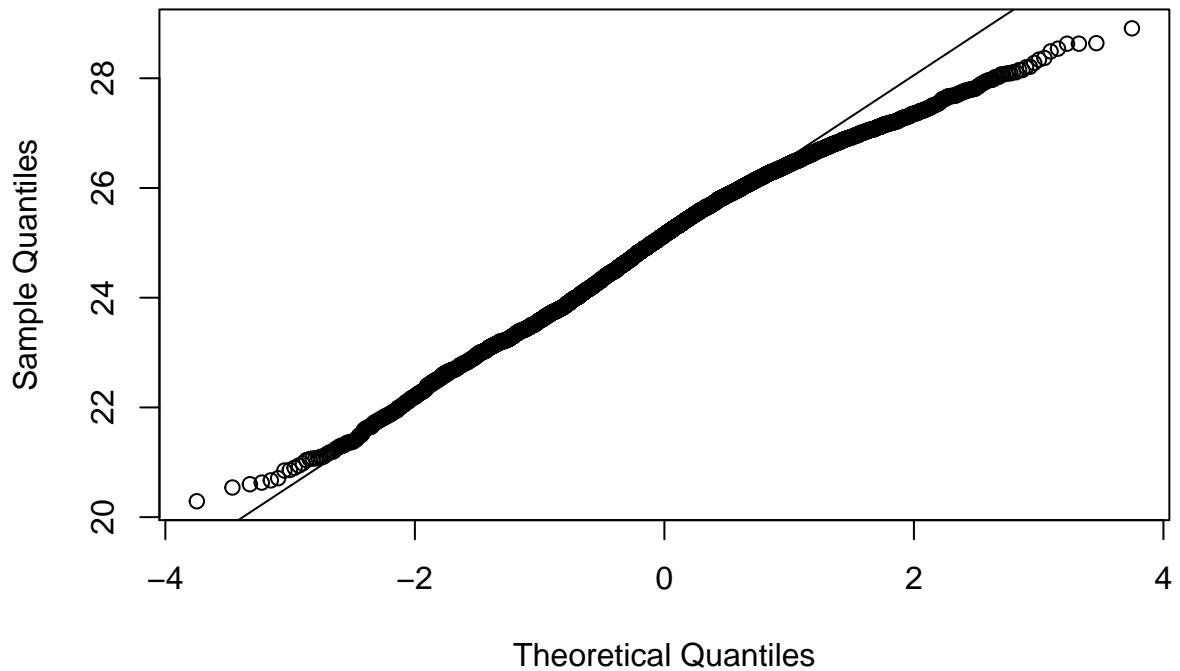
```
ggplot(HawaiiWaterCleanOahu) +
  geom_histogram(aes(x =Temperature), binwidth=0.3)
```



### QQNorm for Temperature

```
qqnorm(HawaiiWaterCleanOahu$Temperature)
qqline(HawaiiWaterCleanOahu$Temperature)
```

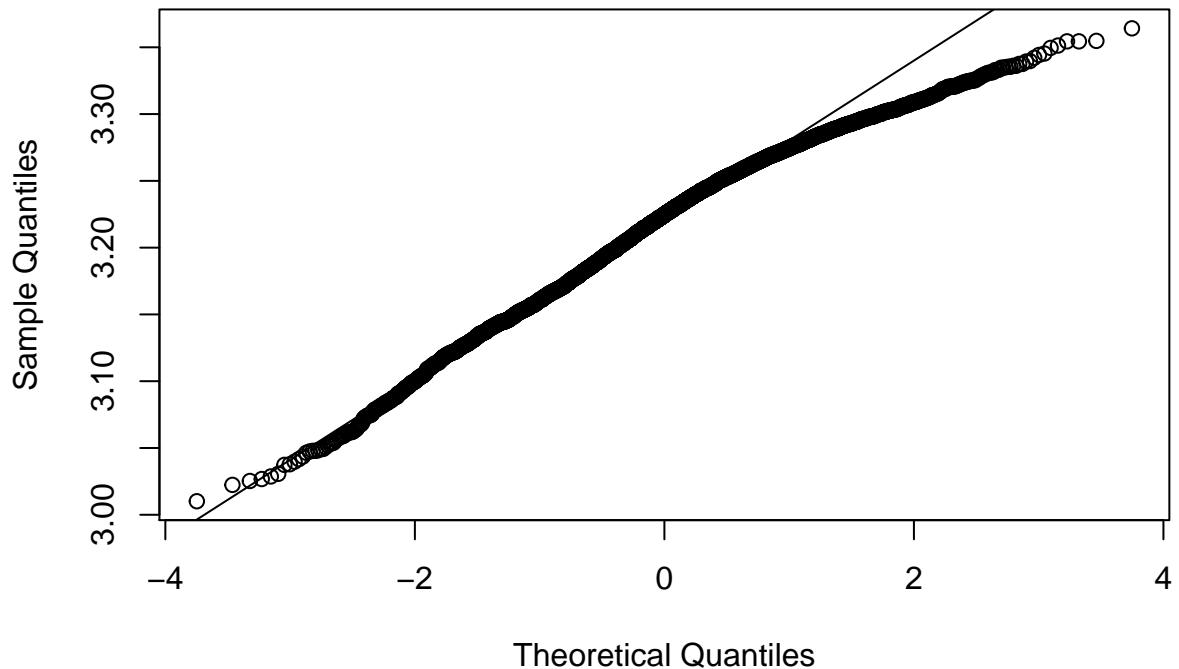
## Normal Q-Q Plot



LogTransform Temperature-doesn't look any better

```
qqnorm(log(HawaiiWaterCleanOahu$Temperature))
qqline(log(HawaiiWaterCleanOahu$Temperature))
```

## Normal Q-Q Plot



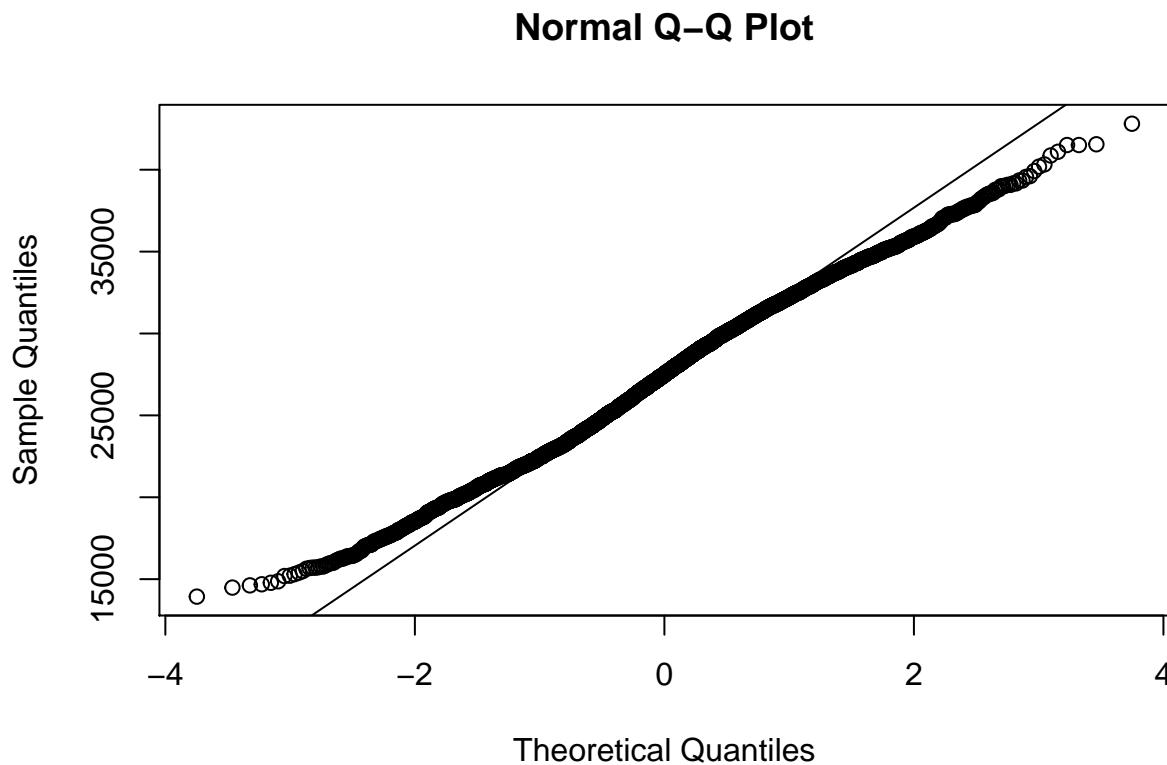
Try using powerTransform function to see what transformation power will make “Temperature” normally distributed

```
summary(powerTransform(HawaiiWaterCleanOahu$Temperature))

## bcPower Transformation to Normality
##                                     Est Power Rounded Pwr Wald Lwr Bnd
## HawaiiWaterCleanOahu$Temperature    3.1789      3.18     2.7406
##                                     Wald Upr Bnd
## HawaiiWaterCleanOahu$Temperature    3.6172
##
## Likelihood ratio test that transformation parameter is equal to 0
##   (log transformation)
##                               LRT df      pval
## LR test, lambda = (0) 208.668  1 < 2.22e-16
##
## Likelihood ratio test that no transformation is needed
##                               LRT df      pval
## LR test, lambda = (1) 97.03683  1 < 2.22e-16
```

Raise Temperature to 3.17 power-looks worse

```
qqnorm((HawaiiWaterCleanOahu$Temperature)^3.17)
qqline((HawaiiWaterCleanOahu$Temperature)^3.17)
```



Perform Shapiro Wilks Normality test for first 5,000 Temperature observations

```
shapiro.test(HawaiiWaterCleanOahu$Temperature[0:5000])
```

```
##
##  Shapiro-Wilk normality test
##
## data: HawaiiWaterCleanOahu$Temperature[0:5000]
## W = 0.98892, p-value < 2.2e-16
```

Summary of Temperature

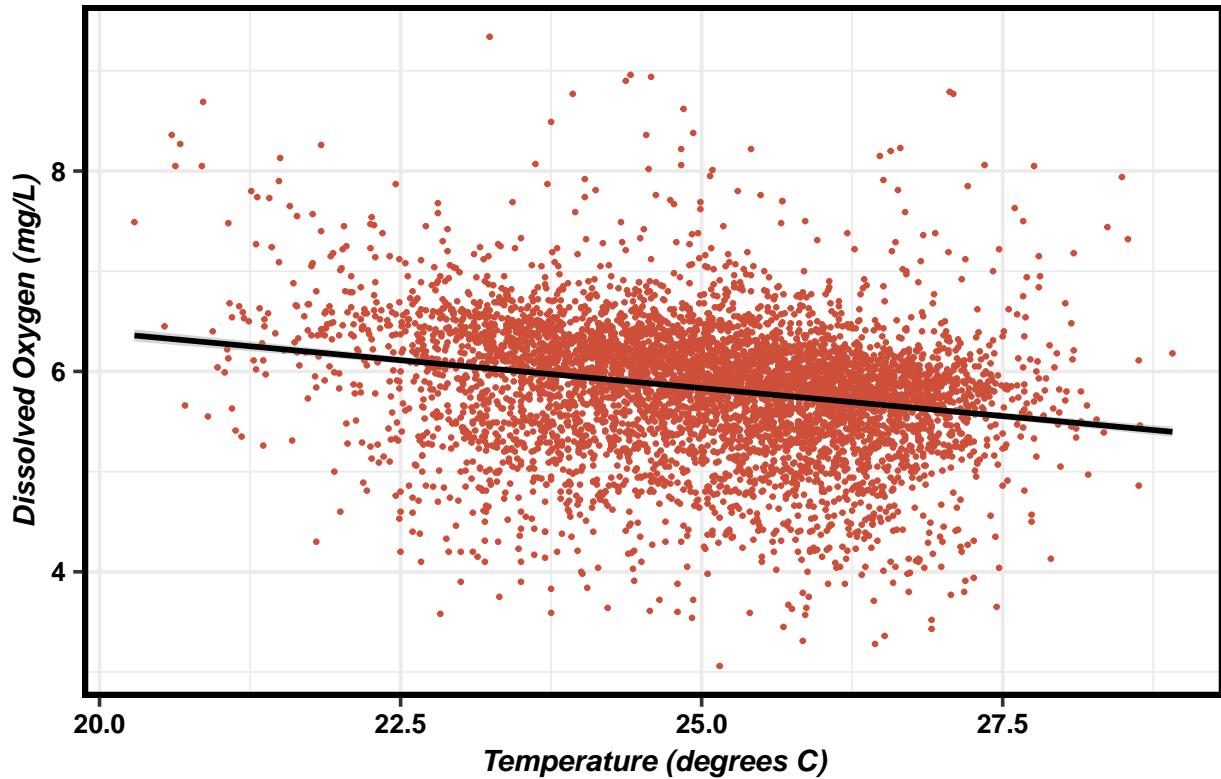
```
summary(HawaiiWaterCleanOahu$Temperature)
```

```
##      Min. 1st Qu. Median    Mean 3rd Qu.    Max.
##    20.29   24.05  25.14  25.04  26.07  28.91
```

## Plot Temperature against DO

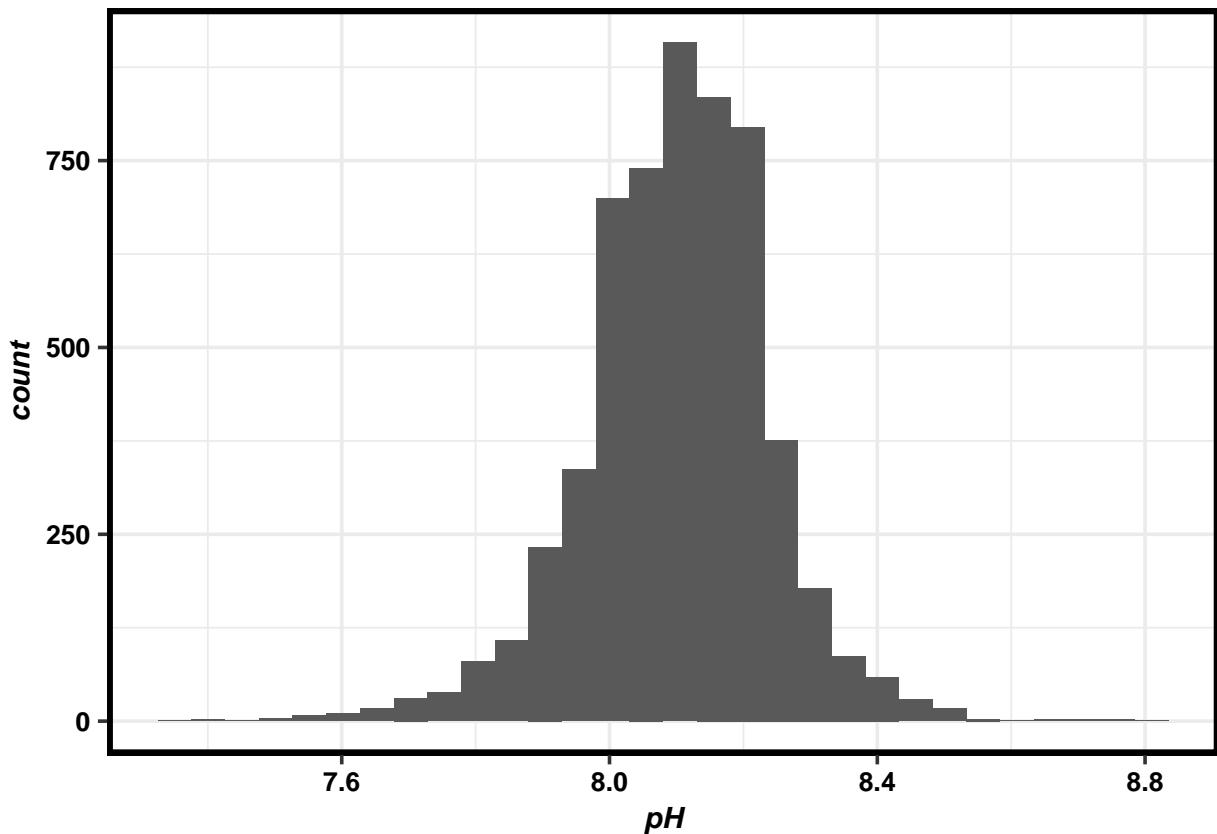
```
TempbyDO <-
  ggplot(HawaiiWaterCleanOahu, aes(x =Temperature, y =DO)) +
  geom_point(color="tomato3", alpha=1, size=0.5) +
  geom_smooth(method=lm, color="black") +
  labs(title="The Effect of Temperature on DO Concentrations across Oahu", x="Temperature (degrees C)",
print(TempbyDO)
```

### The Effect of Temperature on DO Concentrations across Oahu



## pH

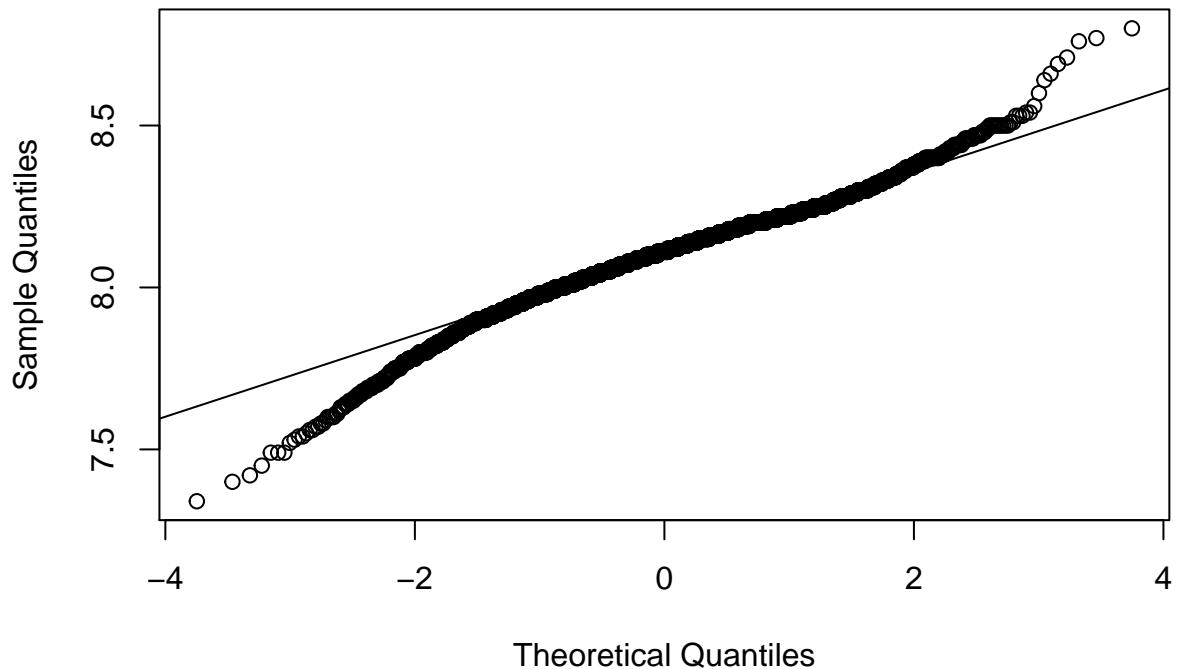
```
ggplot(HawaiiWaterCleanOahu) +
  geom_histogram(aes(x =pH))
```



### QQNorm of pH

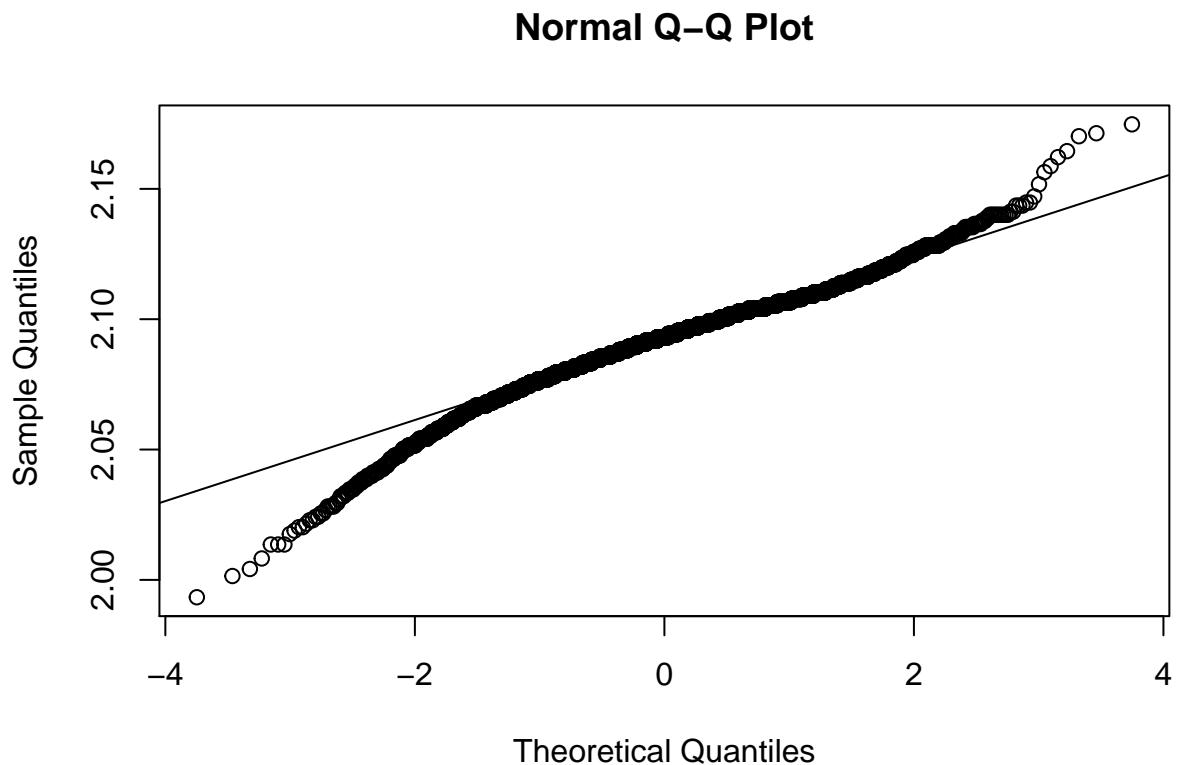
```
qqnorm(HawaiiWaterCleanOahu$pH)
qqline(HawaiiWaterCleanOahu$pH)
```

### Normal Q-Q Plot



Log pH variable-doesn't look any better

```
qqnorm(log(HawaiiWaterCleanOahu$pH))  
qqline(log(HawaiiWaterCleanOahu$pH))
```



#### Shapiro Wilks Test for pH

```
shapiro.test(HawaiiWaterCleanOahu$pH[0:5000])

##
##  Shapiro-Wilk normality test
##
## data: HawaiiWaterCleanOahu$pH[0:5000]
## W = 0.97762, p-value < 2.2e-16
```

#### Summary of pH

```
summary(HawaiiWaterCleanOahu$pH)

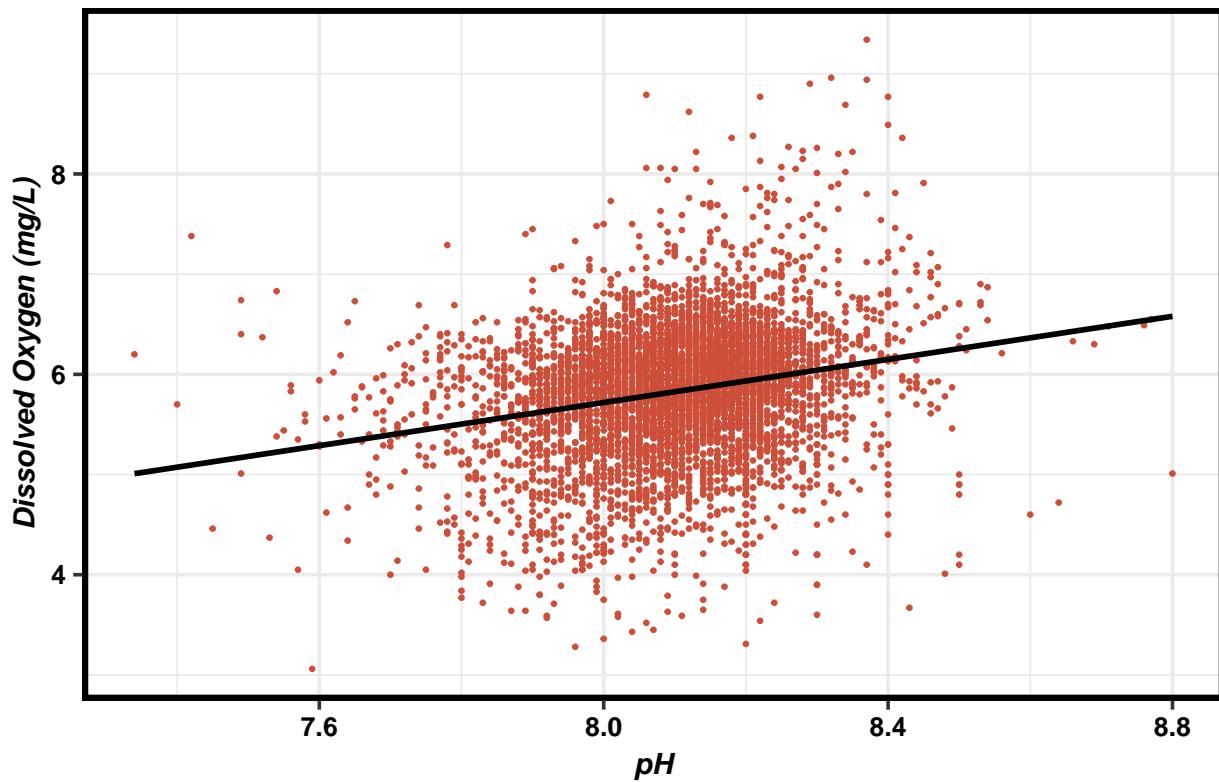
##      Min. 1st Qu. Median     Mean 3rd Qu.    Max.
## 7.340   8.020  8.110  8.103  8.190  8.800
```

#### Plot pH against DO

```
pHbyDO <-
  ggplot(HawaiiWaterCleanOahu, aes(x =pH, y =DO)) +
  geom_point(color="tomato3", alpha=1, size=0.5) +
```

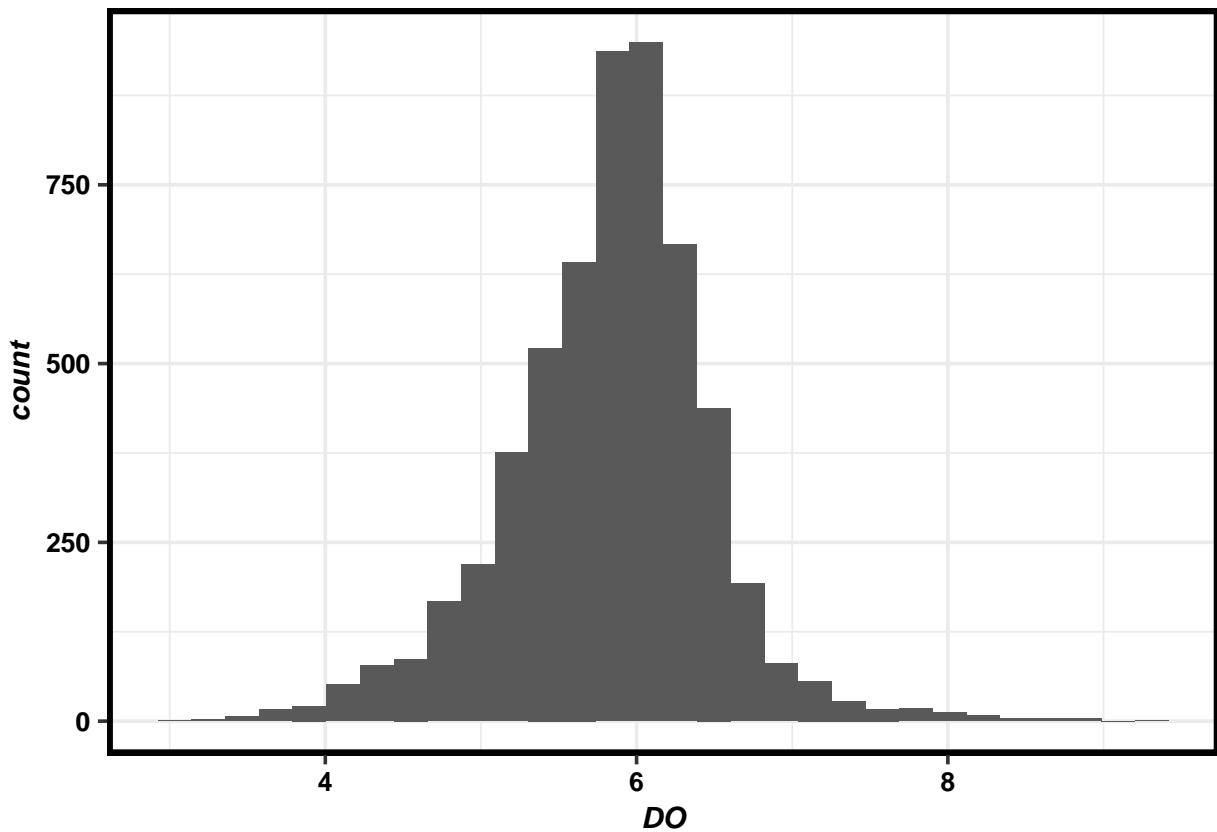
```
geom_smooth(method=lm, color="black", se=FALSE) +
  labs(title="The Effect of pH on DO Concentrations across Oahu", x="pH", y="Dissolved Oxygen (mg/L)")
print(pHbyDO)
```

## The Effect of pH on DO Concentrations across Oahu



## Dissolved Oxygen Concentrations

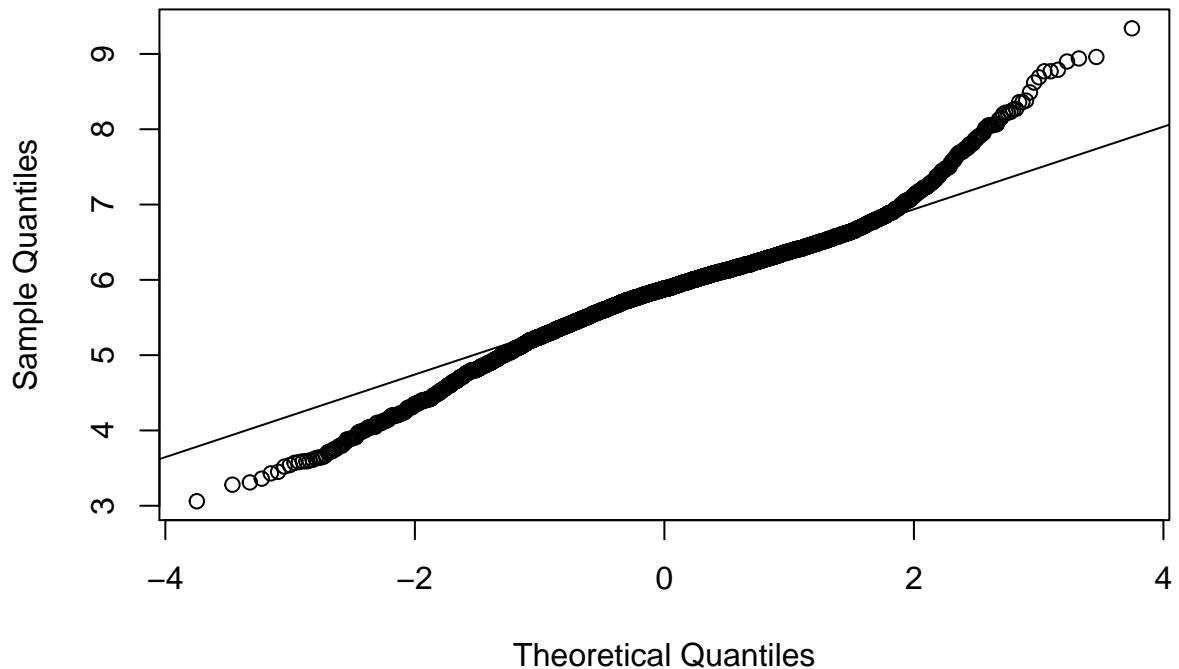
```
ggplot(HawaiiWaterCleanOahu) +
  geom_histogram(aes(x =DO))
```



### QQNorm of DO

```
qqnorm(HawaiiWaterCleanOahu$DO)
qqline(HawaiiWaterCleanOahu$DO)
```

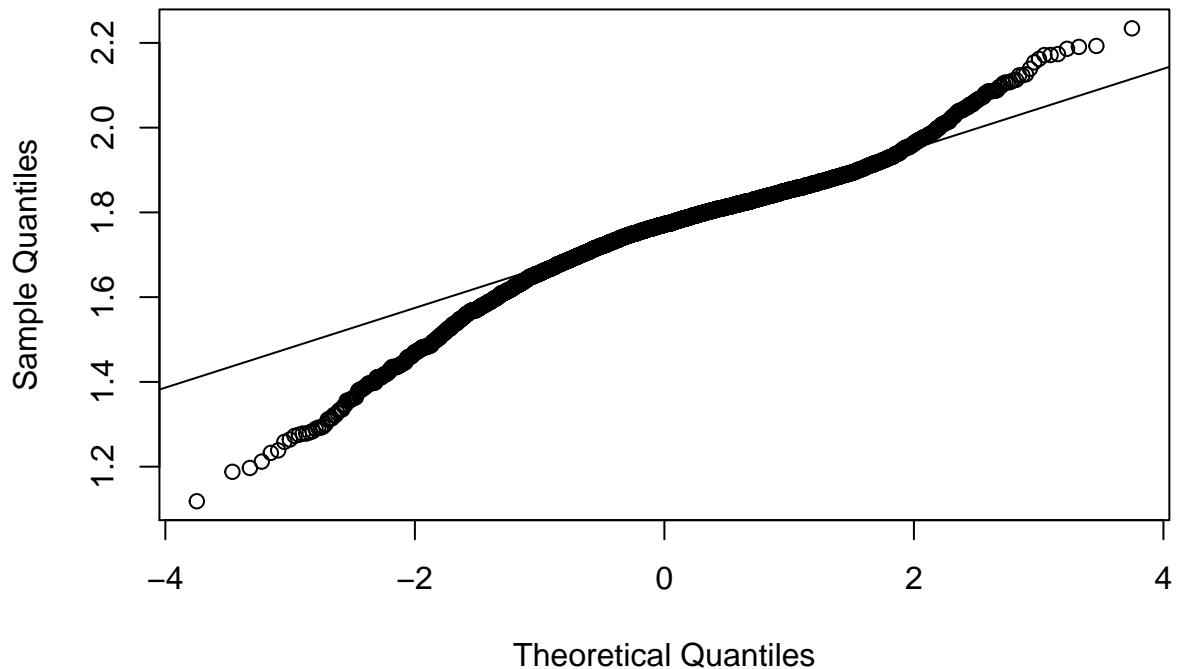
## Normal Q-Q Plot



### Log Transform Dependent Variable DO

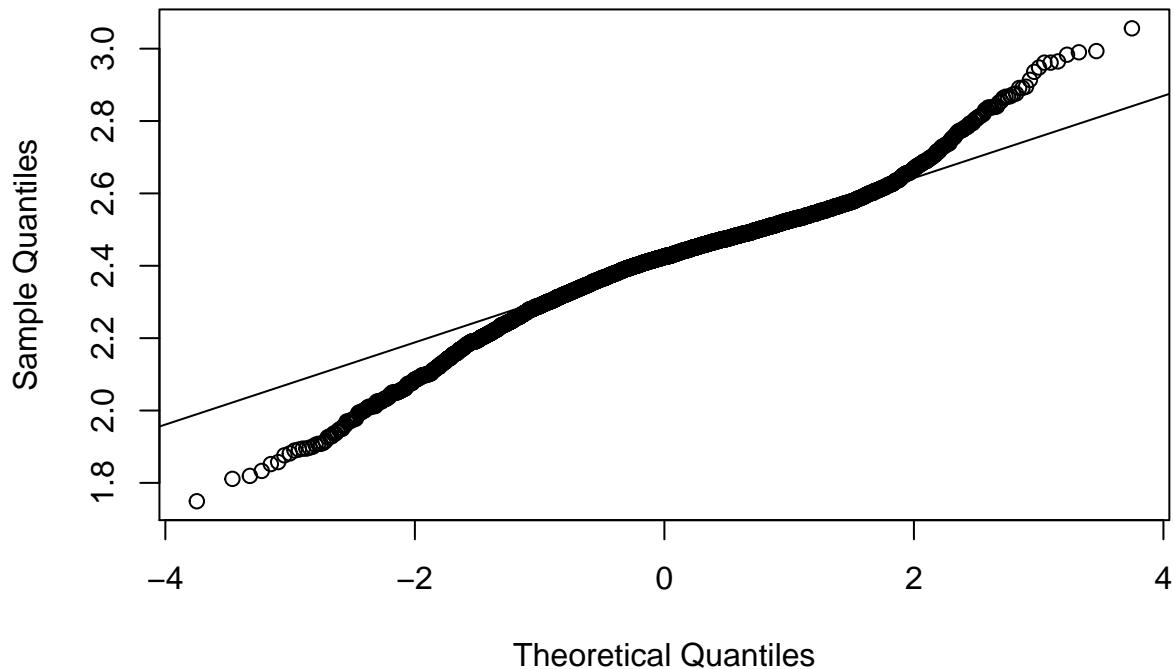
```
qqnorm(log(HawaiiWaterCleanOahu$DO))  
qqline(log(HawaiiWaterCleanOahu$DO))
```

### Normal Q-Q Plot



```
qqnorm(sqrt(HawaiiWaterCleanOahu$DO))  
qqline(sqrt(HawaiiWaterCleanOahu$DO))
```

## Normal Q-Q Plot



Use powertransform Function to generate estimation of Power Lambda that will normalize DV

```
summary(powerTransform(HawaiiWaterCleanOahu$DO))

## bcPower Transformation to Normality
##                               Est Power Rounded Pwr Wald Lwr Bnd Wald Upr Bnd
## HawaiiWaterCleanOahu$DO      1.286      1.29      1.1415      1.4306
##
## Likelihood ratio test that transformation parameter is equal to 0
## (log transformation)
##                               LRT df      pval
## LR test, lambda = (0) 305.8359  1 < 2.22e-16
##
## Likelihood ratio test that no transformation is needed
##                               LRT df      pval
## LR test, lambda = (1) 15.05389  1 0.00010448
```

Shapiro Wilks Test for DO

```
shapiro.test(HawaiiWaterCleanOahu$DO[0:5000])

##
```

```
## Shapiro-Wilk normality test
##
## data: HawaiiWaterCleanOahu$DO[0:5000]
## W = 0.97535, p-value < 2.2e-16
```

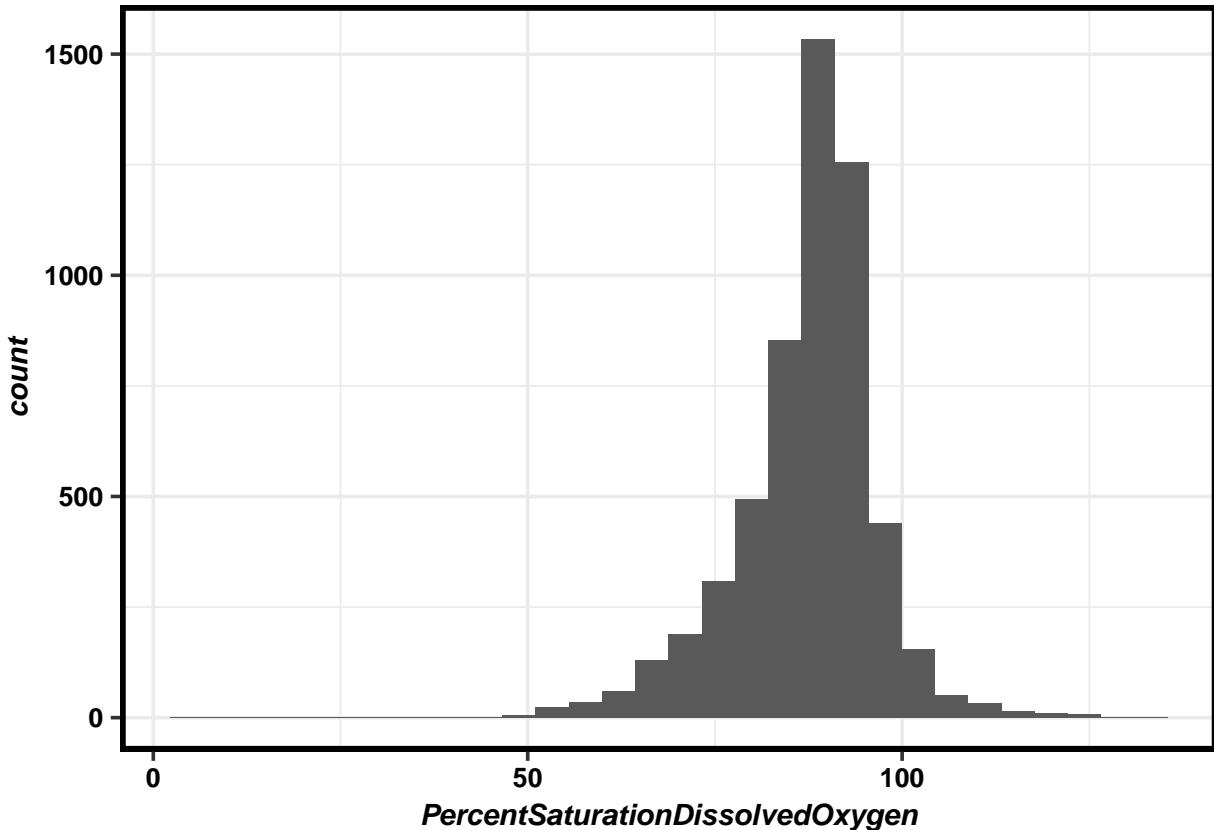
## Summary of DO

```
summary(HawaiiWaterCleanOahu$DO)
```

```
##      Min.   1st Qu.    Median     Mean   3rd Qu.   Max.
##      3.060   5.470   5.880   5.829   6.210   9.340
```

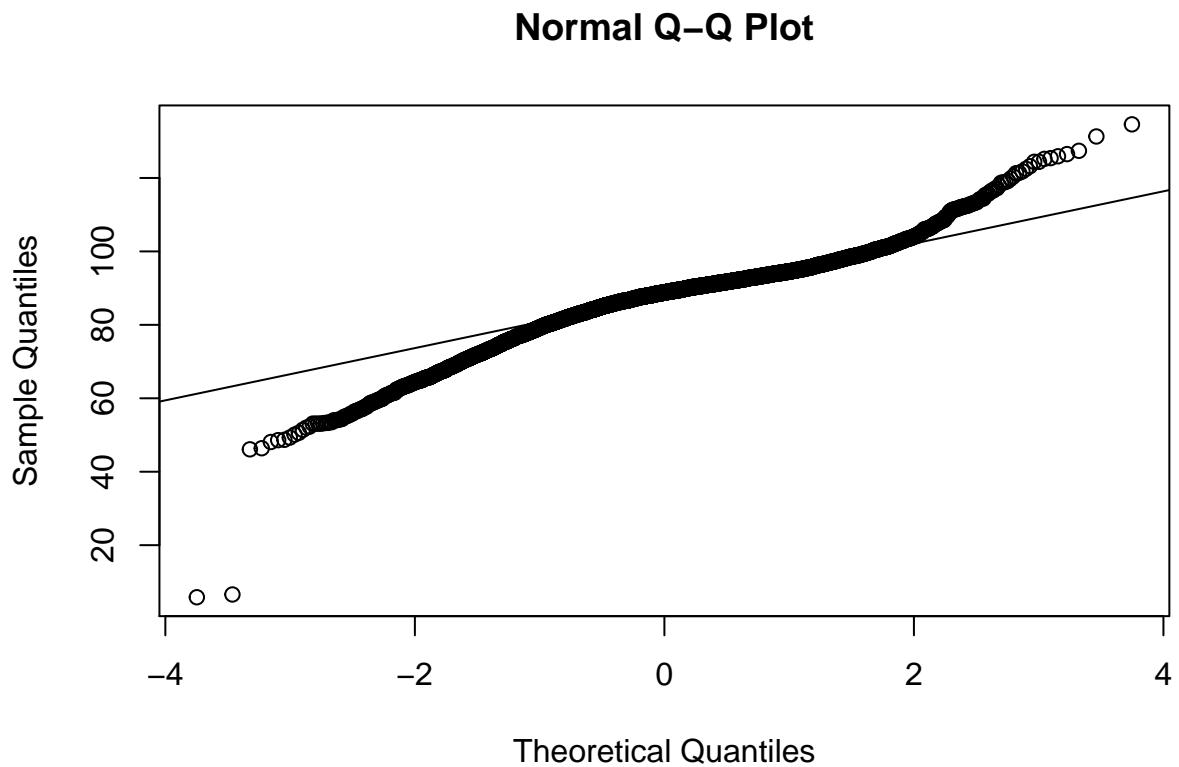
## Percent Saturation of Dissolved Oxygen

```
ggplot(HawaiiWaterCleanOahu) +
  geom_histogram(aes(x = PercentSaturationDissolvedOxygen))
```



## QQnorm of Percent Saturation of Dissolved Oxygen

```
qqnorm(HawaiiWaterCleanOahu$PercentSaturationDissolvedOxygen)
qqline(HawaiiWaterCleanOahu$PercentSaturationDissolvedOxygen)
```



### Shapiro Test for Percent Saturation Dissolved Oxygen

```
shapiro.test(HawaiiWaterCleanOahu$PercentSaturationDissolvedOxygen[0:5000])

##
##  Shapiro-Wilk normality test
##
## data: HawaiiWaterCleanOahu$PercentSaturationDissolvedOxygen[0:5000]
## W = 0.95108, p-value < 2.2e-16
```

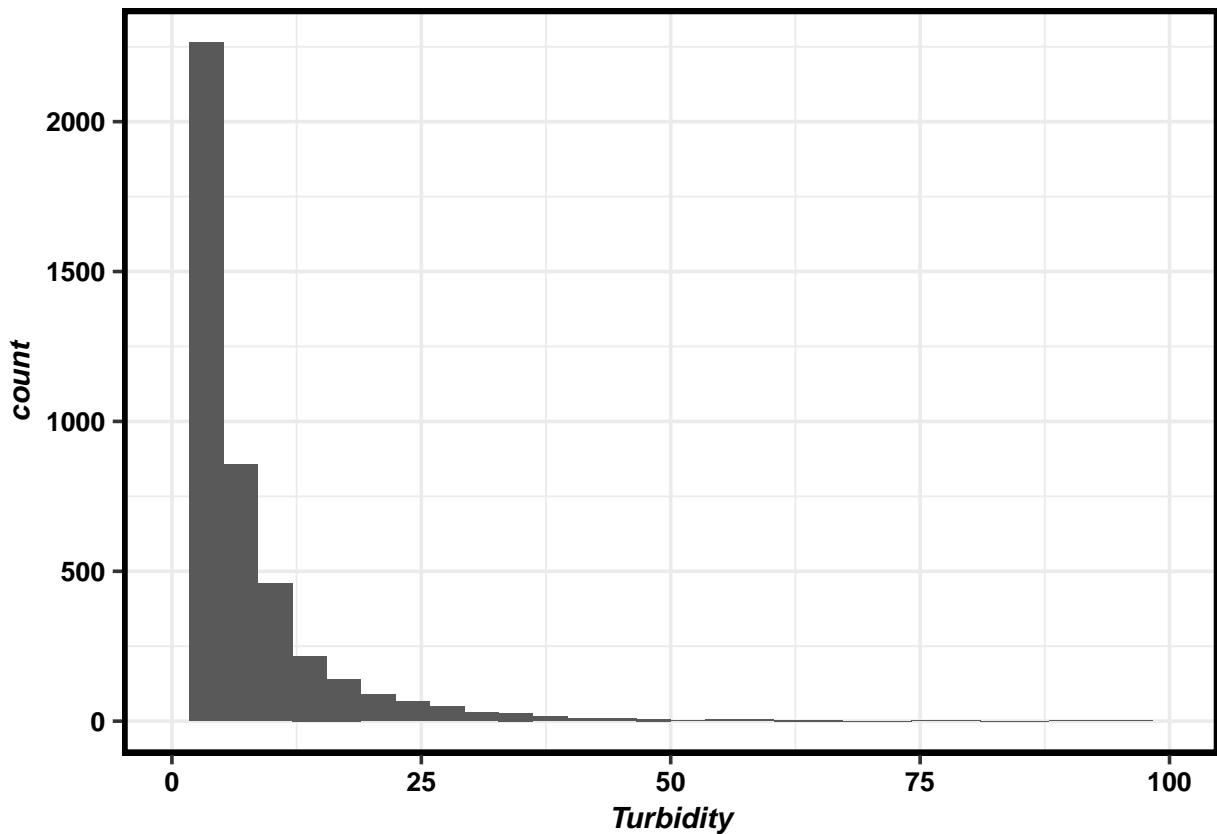
### Summary of Percent Saturation of Dissolved Oxygen

```
summary(HawaiiWaterCleanOahu$PercentSaturationDissolvedOxygen)

##      Min. 1st Qu. Median     Mean 3rd Qu.    Max.
##      5.82   83.10  88.80   87.38  92.70 134.60
```

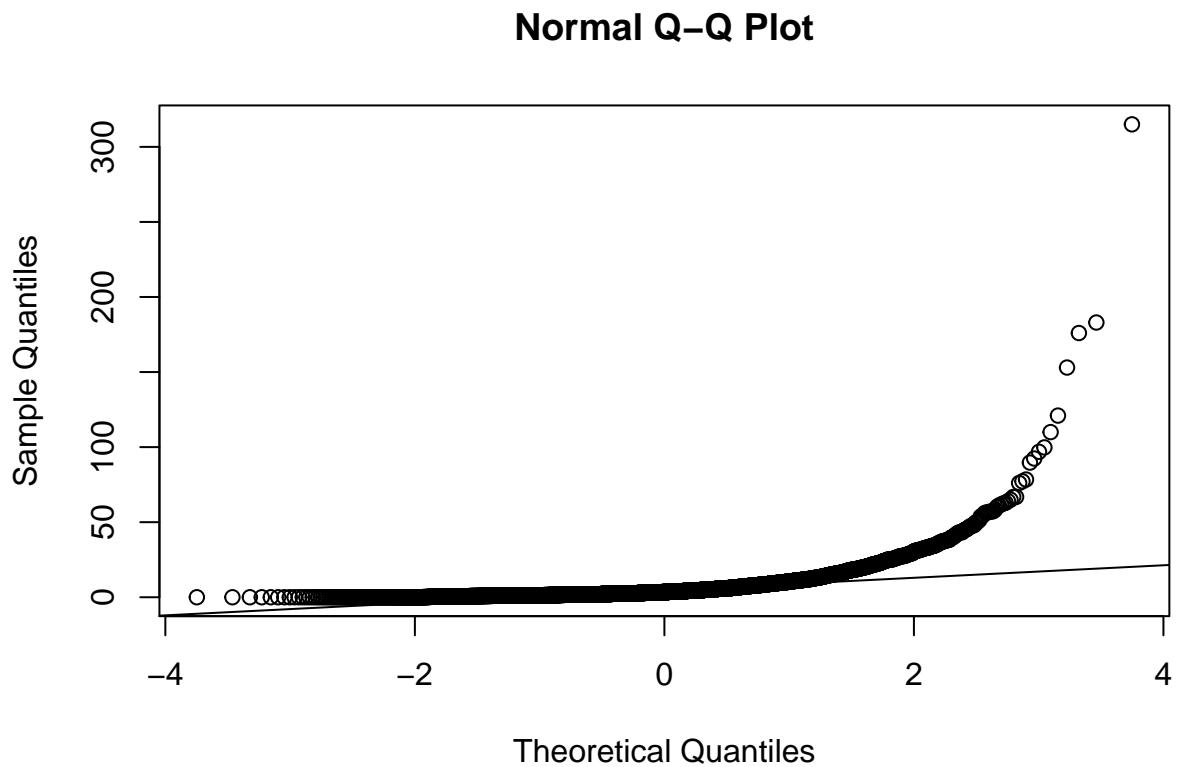
### Turbidity

```
ggplot(HawaiiWaterCleanOahu) +
  geom_histogram(aes(x = Turbidity)) +
  scale_x_continuous(limits = c(0, 100))
```



### QQNorm of Turbidity

```
qqnorm(HawaiiWaterCleanOahu$Turbidity)
qqline(HawaiiWaterCleanOahu$Turbidity)
```



#### Shapiro Test for Turbidity

```
shapiro.test(HawaiiWaterCleanOahu$Turbidity[0:5000])

##
##  Shapiro-Wilk normality test
##
## data: HawaiiWaterCleanOahu$Turbidity[0:5000]
## W = 0.55091, p-value < 2.2e-16
```

#### Summary of Turbidity

```
summary(HawaiiWaterCleanOahu$Turbidity)

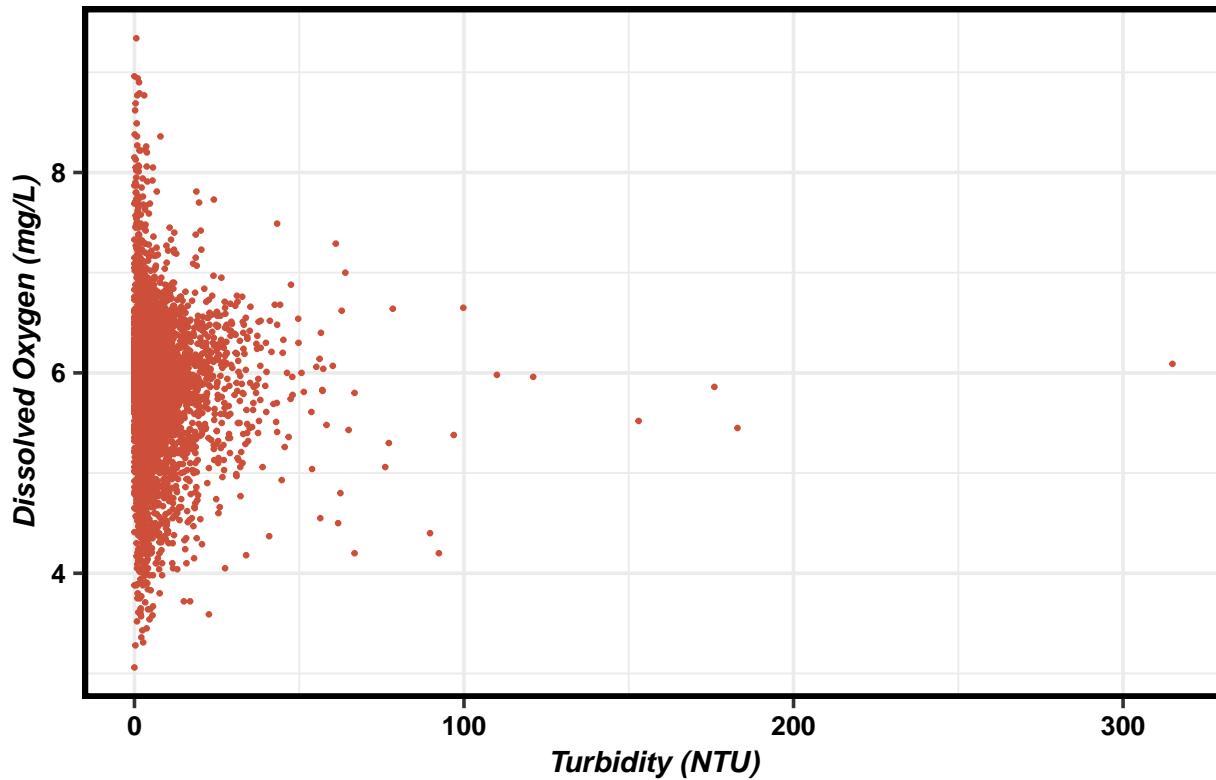
##      Min. 1st Qu. Median     Mean 3rd Qu.    Max.
## 0.000   1.800   3.440   6.331   7.425 315.000
```

#### Plot Turbidity against DO

```
TurbiditybyDO <-
ggplot(HawaiiWaterCleanOahu, aes(x =Turbidity, y =DO)) +
  geom_point(color="tomato3", alpha=1, size=0.5)
```

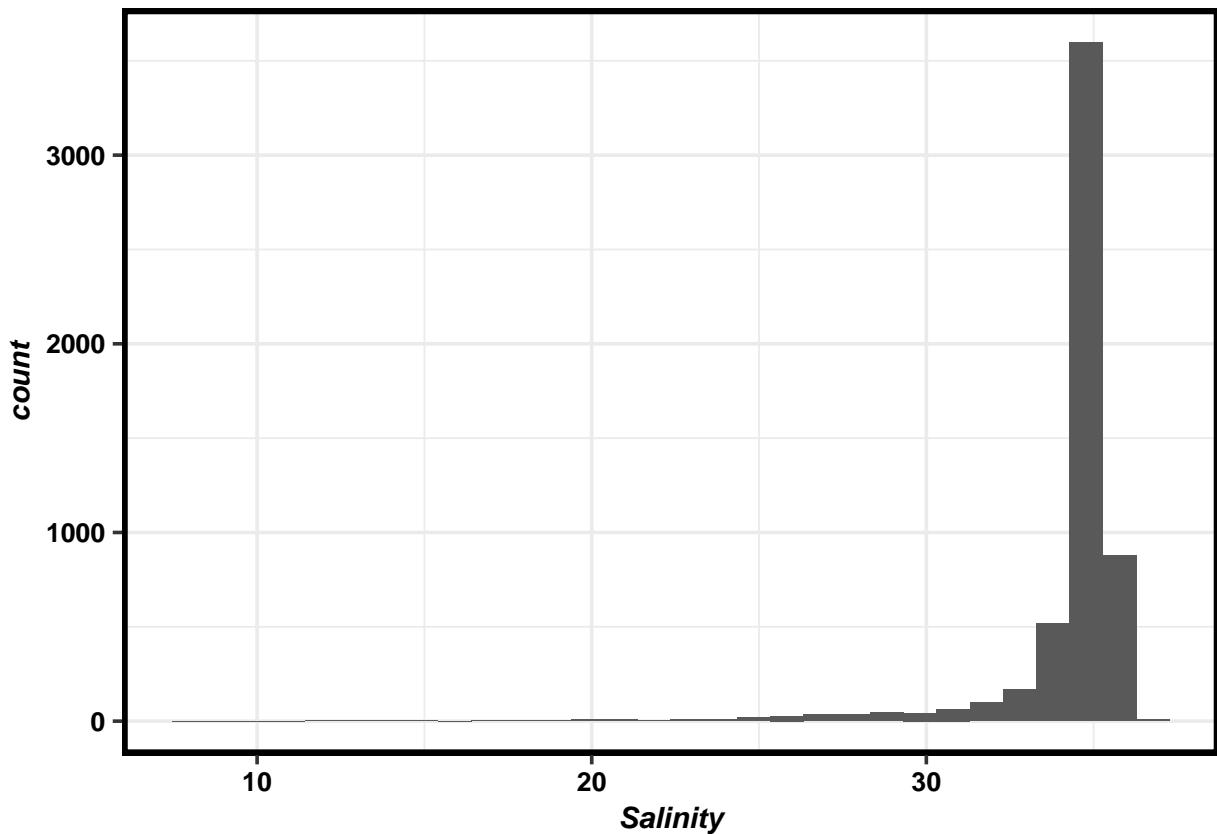
```
  labs(title="The Effect of Turbidity on DO Concentrations across Oahu", x="Turbidity (NTU)", y="Dissolved Oxygen (mg/L)")  
print(TurbiditybyDO)
```

## The Effect of Turbidity on DO Concentrations across Oahu



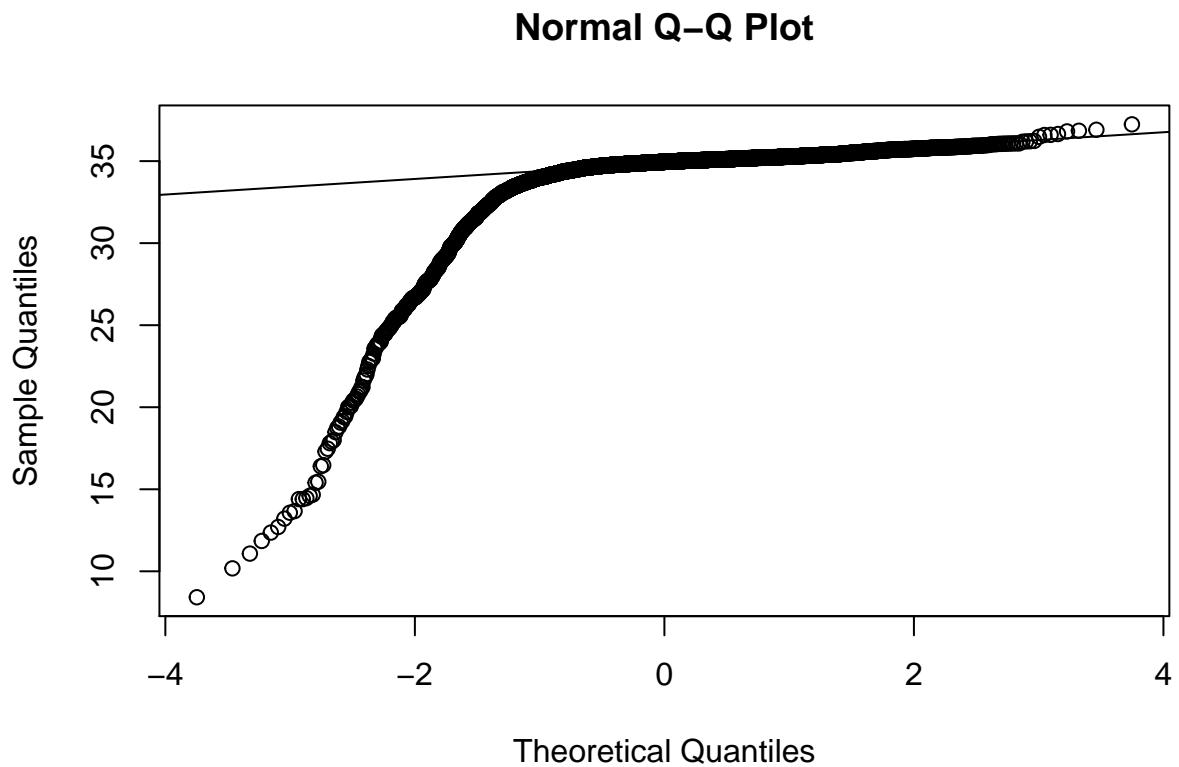
## Salinity

```
ggplot(HawaiiWaterCleanOahu) +  
  geom_histogram(aes(x = Salinity))
```



### QQNorm of Salinity

```
qqnorm(HawaiiWaterCleanOahu$Salinity)
qqline(HawaiiWaterCleanOahu$Salinity)
```



#### Shapiro Test for Salinity

```
shapiro.test(HawaiiWaterCleanOahu$Salinity[0:5000])

##
##  Shapiro-Wilk normality test
##
## data: HawaiiWaterCleanOahu$Salinity[0:5000]
## W = 0.42958, p-value < 2.2e-16
```

#### Summary of Salinity

```
summary(HawaiiWaterCleanOahu$Salinity)

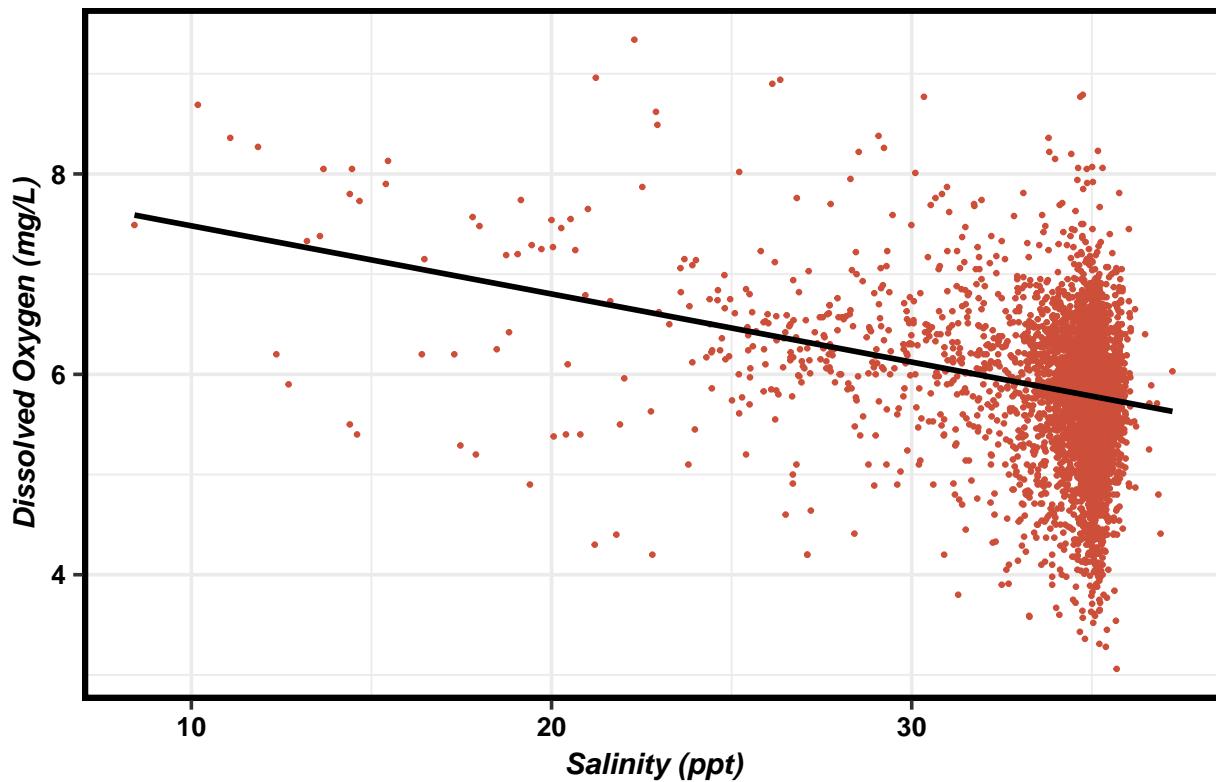
##      Min. 1st Qu. Median     Mean 3rd Qu.    Max.
##     8.42   34.54  34.95   34.31  35.18   37.24
```

#### Plot Salinity Against DO

```
SalinitybyDO <-
  ggplot(HawaiiWaterCleanOahu, aes(x = Salinity, y = DO)) +
  geom_point(color="tomato3", alpha=1, size=0.5)
```

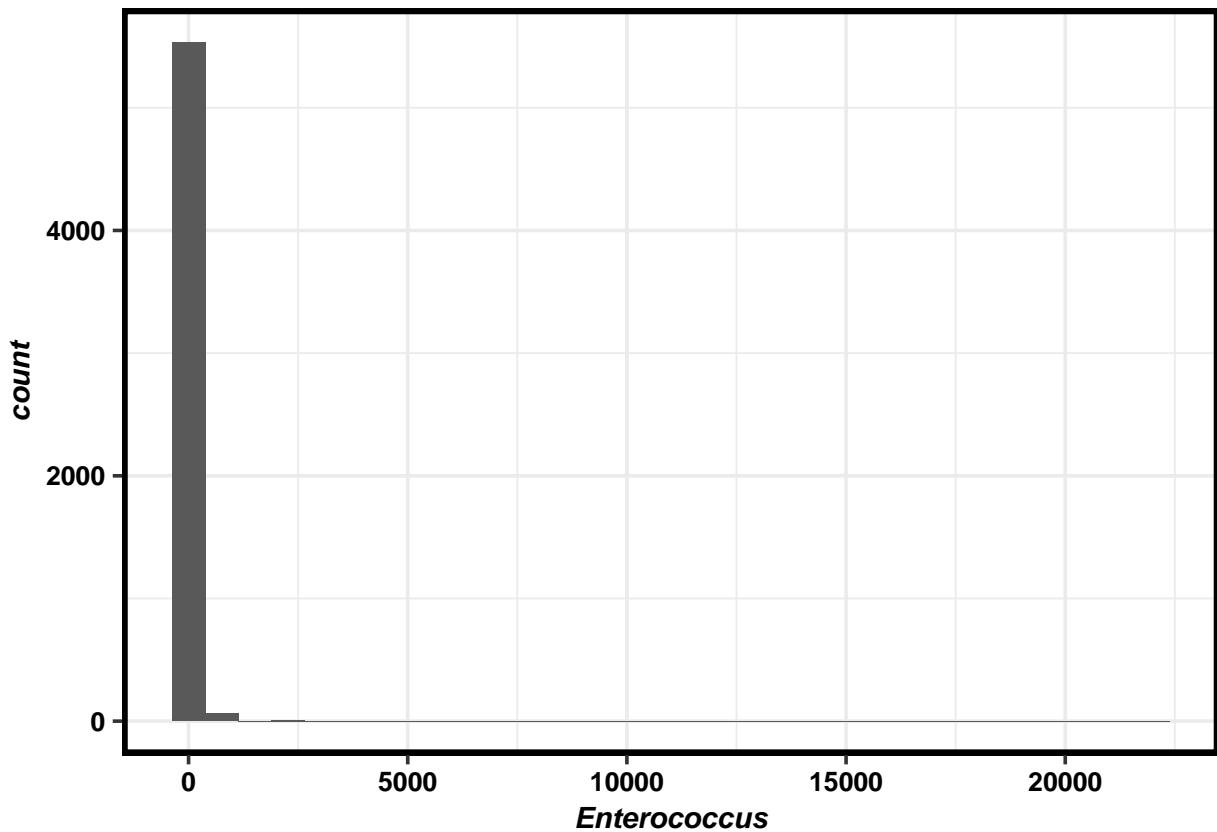
```
geom_smooth(method=lm, color="black", se=FALSE) +
  labs(title="The Effect of Salinity on DO Concentrations across Oahu", x="Salinity (ppt)", y="Dissolved Oxygen (mg/L)")
```

## The Effect of Salinity on DO Concentrations across Oahu



## Enterococcus

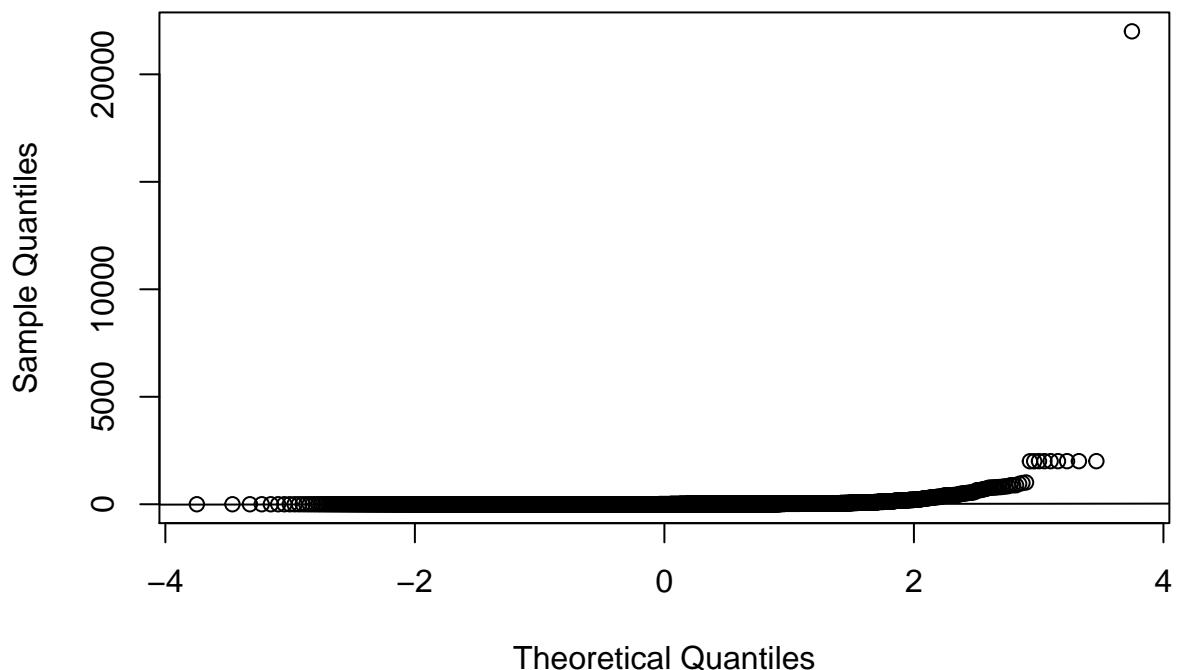
```
ggplot(HawaiiWaterCleanOahu) +
  geom_histogram(aes(x = Enterococcus))
```



### QQNorm of Enterococcus

```
qqnorm(HawaiiWaterCleanOahu$Enterococcus)
qqline(HawaiiWaterCleanOahu$Enterococcus)
```

## Normal Q-Q Plot



## Shapiro Test for Enterococcus

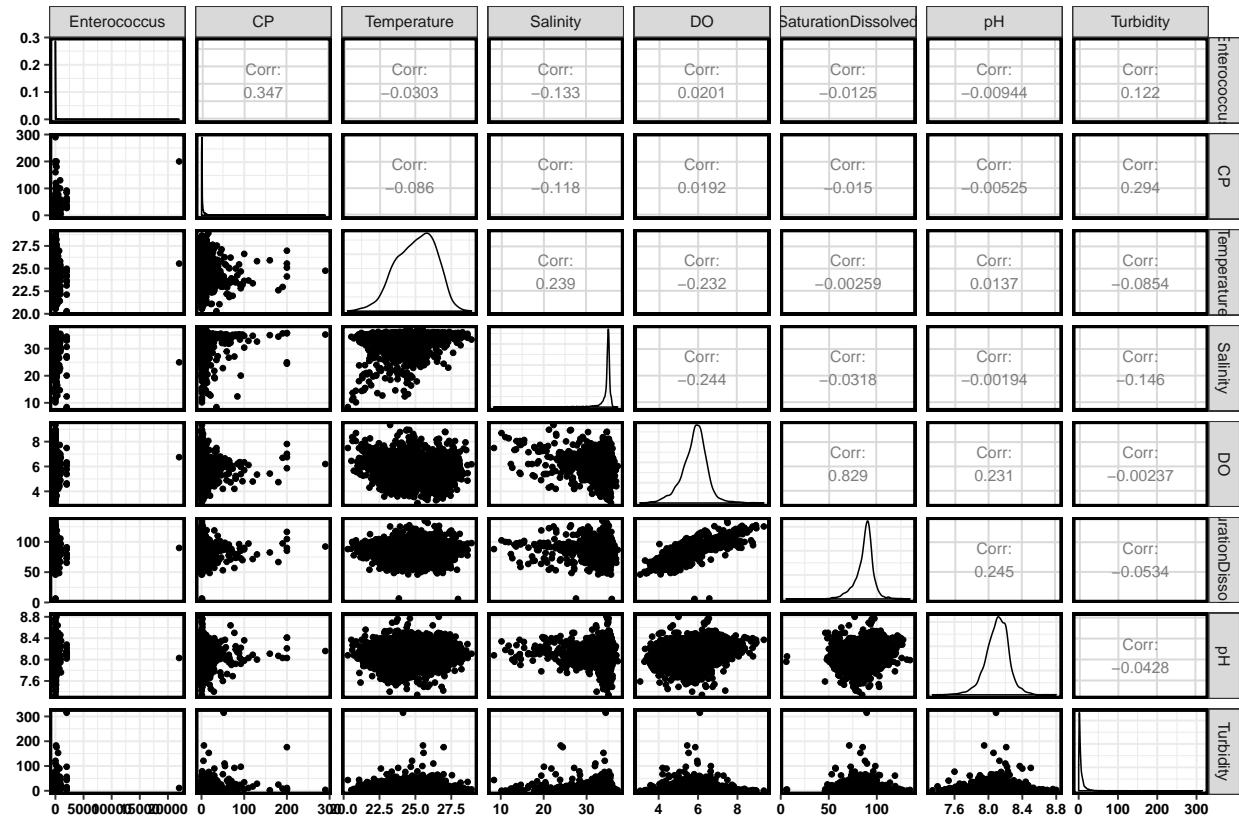
```
shapiro.test(HawaiiWaterCleanOahu$Enterococcus[0:5000])  
  
##  
## Shapiro-Wilk normality test  
##  
## data: HawaiiWaterCleanOahu$Enterococcus[0:5000]  
## W = 0.03267, p-value < 2.2e-16
```

## Summary of Enterococcus

```
summary(HawaiiWaterCleanOahu$Enterococcus)  
  
##      Min.   1st Qu.    Median     Mean   3rd Qu.    Max.  
##      0.30    2.30    3.30   28.51   10.00 22000.00
```

## Correlation Plot of Data

```
ggpairs(HawaiiWaterCleanOahu, columns = 8:15)
```



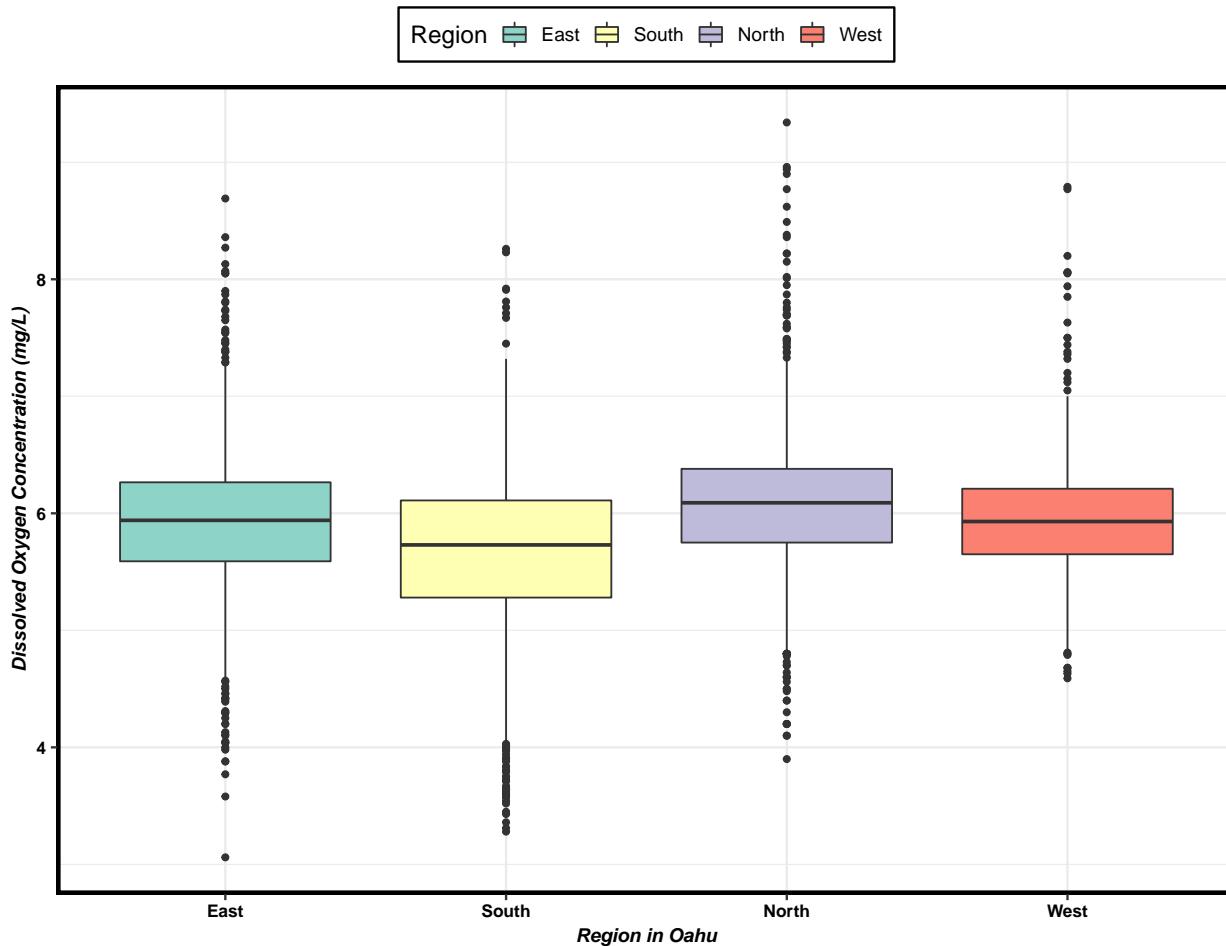
###Enterococcus and DO have almost no correlation: 0.0201 ###CP and DO have almost no correlation: 0.0192 ###Temperature and DO have a low negative correlation: -0.232 ###Salinity and DO have a low negative correlation: -0.244 ###pH and DO have a low positive correlation: 0.231 ###Turbidity and DO have almost no correlation: -0.002 ###Percent Saturation DO and DO are multicollineated because they are related variables, so I won't be including Percent Saturation of DO in my analysis

## Exploratory Boxplot showing Range of DO Concentrations by Region

```
D0Boxplot<-ggplot(HawaiiWaterCleanOahu) +
  geom_boxplot(aes(x=Region, y=DO, fill=Region)) +
  labs(title="Effect of Region on Range of Dissolved Oxygen Concentrations in Oahu", x="Region in Oahu",
       theme(legend.title = element_text(colour="IndianRed", size=16, face="bold")) +gabytheme +
       scale_fill_brewer(palette="Set3"))

print(D0Boxplot)
```

### Effect of Region on Range of Dissolved Oxygen Concentrations in Oahu



### Research Question Number 1:

Which of the parameters have a relationship with dissolved oxygen concentrations? Of these relevant parameters, which have the most significant effect on dissolved oxygen concentrations over time?

When I did the full maximal model including ALL interactions, the AIC of the maximal model was higher than the 12 subsequent reduced models—>Too many parameters with all of the interactions, so I decided to not include interactions.

#### Full Maximal Model

```
attach(HawaiiWaterCleanOahu)
HawaiiMod<-glm(DO~Enterococcus + Temperature + Salinity + pH + Turbidity + CP, data=HawaiiWaterCleanOahu,
summary(HawaiiMod)

##
## Call:
## glm(formula = DO ~ Enterococcus + Temperature + Salinity + pH +
##     Turbidity + CP, family = "gaussian", data = HawaiiWaterCleanOahu)
```

```

## 
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.6668 -0.3024  0.0875  0.3569  3.2043
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.352e+00 4.968e-01  2.721  0.00653 **
## Enterococcus -7.018e-06 2.759e-05 -0.254  0.79925
## Temperature -9.138e-02 6.128e-03 -14.912 < 2e-16 ***
## Salinity    -5.703e-02 3.602e-03 -15.834 < 2e-16 ***
## pH           1.078e+00 5.745e-02 18.774 < 2e-16 ***
## Turbidity   -2.328e-03 8.498e-04 -2.740  0.00617 **
## CP          -5.034e-04 8.308e-04 -0.606  0.54454
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for gaussian family taken to be 0.3632317)
## 
## Null deviance: 2384.6 on 5603 degrees of freedom
## Residual deviance: 2033.0 on 5597 degrees of freedom
## AIC: 10237
## 
## Number of Fisher Scoring iterations: 2

```

## Remove Enterococcus Parameter

```
HawaiiMod2<-update(HawaiiMod,.~.-Enterococcus)
summary(HawaiiMod2)
```

```

## 
## Call:
## glm(formula = DO ~ Temperature + Salinity + pH + Turbidity +
##       CP, family = "gaussian", data = HawaiiWaterCleanOahu)
## 
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.6661 -0.3022  0.0873  0.3564  3.2043
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.3487262 0.4966506  2.716  0.00664 **
## Temperature -0.0914171 0.0061257 -14.924 < 2e-16 ***
## Salinity    -0.0569437 0.0035840 -15.888 < 2e-16 ***
## pH           1.0786085 0.0574404 18.778 < 2e-16 ***
## Turbidity   -0.0023307 0.0008497 -2.743  0.00611 **
## CP          -0.0005717 0.0007862 -0.727  0.46721
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for gaussian family taken to be 0.3631711)
## 
## Null deviance: 2384.6 on 5603 degrees of freedom
```

```

## Residual deviance: 2033.0 on 5598 degrees of freedom
## AIC: 10235
##
## Number of Fisher Scoring iterations: 2

```

## Remove CP Parameter

```

HawaiiMod3<-update(HawaiiMod2,.~.-CP)
summary(HawaiiMod3)

##
## Call:
## glm(formula = DO ~ Temperature + Salinity + pH + Turbidity, family = "gaussian",
##      data = HawaiiWaterCleanOahu)
##
## Deviance Residuals:
##       Min     1Q   Median     3Q    Max
## -2.6636 -0.3021  0.0880  0.3571  3.2042
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.3394910  0.4964673  2.698  0.00700 **
## Temperature -0.0912068  0.0061186 -14.907 < 2e-16 ***
## Salinity    -0.0567715  0.0035760 -15.876 < 2e-16 ***
## pH          1.0782888  0.0574363  18.774 < 2e-16 ***
## Turbidity   -0.0025036  0.0008157 -3.069  0.00215 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.3631405)
##
## Null deviance: 2384.6 on 5603 degrees of freedom
## Residual deviance: 2033.2 on 5599 degrees of freedom
## AIC: 10234
##
## Number of Fisher Scoring iterations: 2

```

## Analysis of Final Model:

A one degree Celsius increase in temperature will result in a decrease in dissolved oxygen concentrations by 0.09 mg/L A one unit increase in Salinity will result in a decrease in dissolved oxygen concentrations by 0.056 mg/L A one unit increase in pH will result in an increase in dissolved oxygen concentrations by 1.08 mg/L A one unit increase in Turbidity will result in a decrease in dissolved oxygen concentrations by 0.0025 mg/L.

## Equation:

Dissolved Oxygen=1.34 -0.09(Temperature)-0.06(Salinity) + 1.08(pH) -0.003(Turbidity) + E (this could be added in figure as a geom\_text)

## AIC Test of all models

```
AIC(HawaiiMod, HawaiiMod2, HawaiiMod3)
```

```
##          df      AIC
## HawaiiMod   8 10237.21
## HawaiiMod2   7 10235.27
## HawaiiMod3   6 10233.80
```

## Partial F-test of all Models

```
anova(HawaiiMod, HawaiiMod2, HawaiiMod3)
```

```
## Analysis of Deviance Table
##
## Model 1: DO ~ Enterococcus + Temperature + Salinity + pH + Turbidity +
##           CP
## Model 2: DO ~ Temperature + Salinity + pH + Turbidity + CP
## Model 3: DO ~ Temperature + Salinity + pH + Turbidity
##   Resid. Df Resid. Dev Df  Deviance
## 1      5597    2033.0
## 2      5598  2033.0 -1 -0.023496
## 3      5599  2033.2 -1 -0.191988
```

## Check for Multicollinearity of Final Model

```
vif(HawaiiMod3)
```

```
## Temperature     Salinity         pH     Turbidity
## 1.063851     1.079097     1.002056     1.026523
```

Temperature, Salinity, pH, and Turbidity have a significant effect on DO Concentrations and are not

## Check for Overdispersion in Final Model

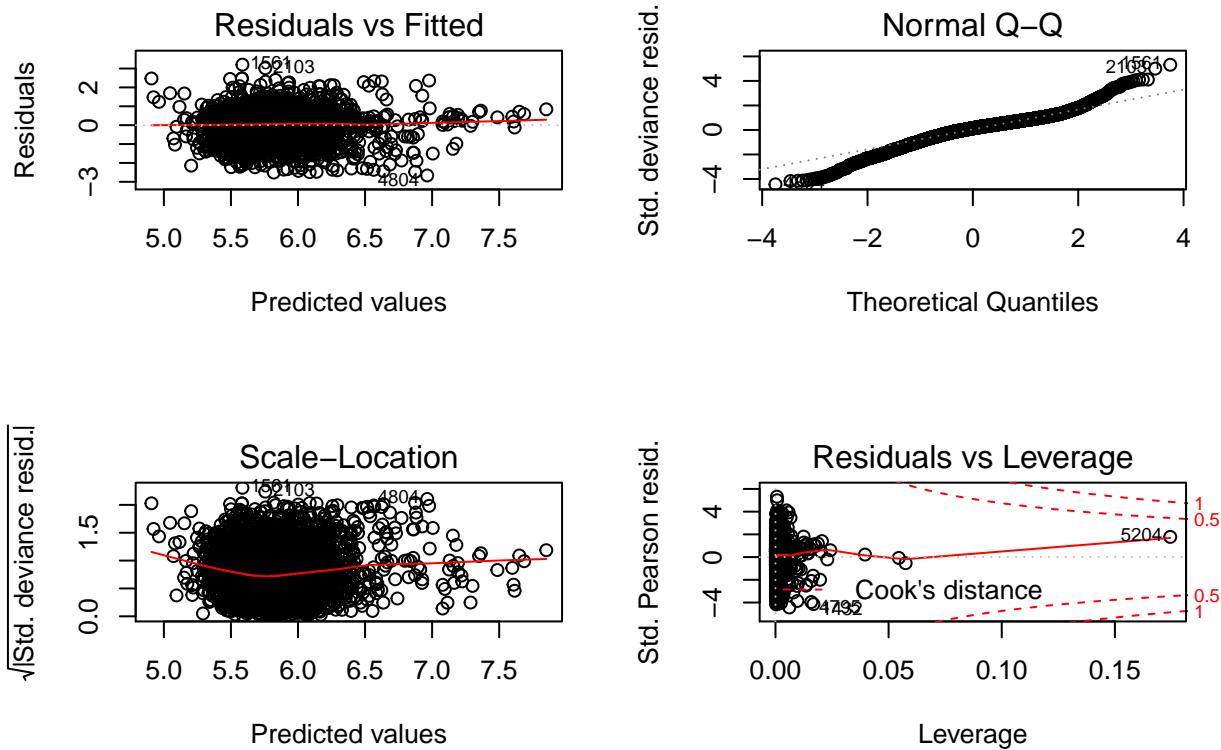
```
HawaiiMod3$deviance/HawaiiMod3$df.residual
```

```
## [1] 0.3631405
```

Final model is not overdispersed because the resulting value is below 1.5

## Check Residuals of HawaiiMod3

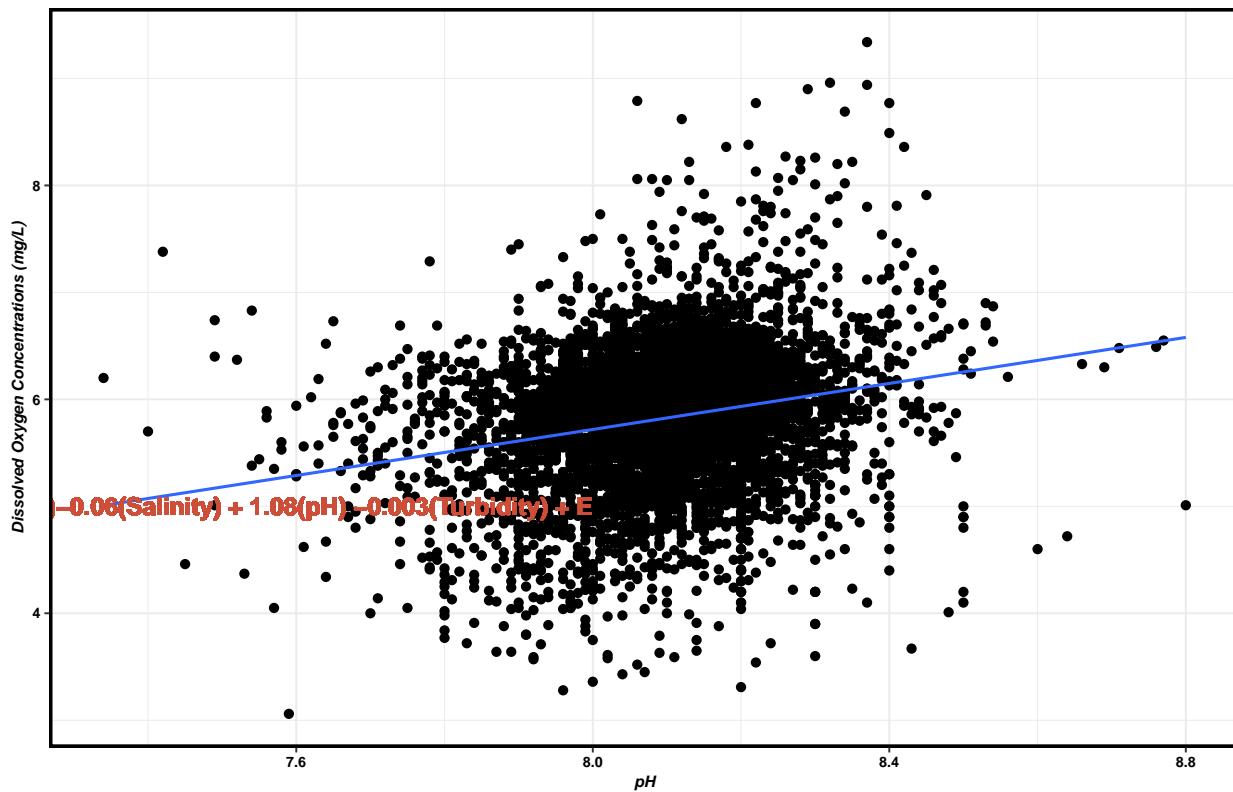
```
par(mfrow=c(2,2))
plot(HawaiiMod3)
```



# Graph

```
pHbyDOPlot<-ggplot(HawaiiWaterCleanOahu, aes(x =pH, y =DO)) +
  geom_point(aes(x =pH, y =DO), size=3, shape=16, alpha=1) +
  geom_smooth(aes(x =pH, y =DO, span=0.1),
              method="lm", se=FALSE, linetype=1, size=1) +
  labs(title="The Effect of pH on DO Concentrations", x="pH",
       y="Dissolved Oxygen Concentrations (mg/L)") +
  geom_text(x =8, y = 5, label = "Dissolved Oxygen=1.34 -0.09(Temperature)-0.06(Salinity) + 1.08(pH) -0
## Warning: Ignoring unknown aesthetics: span
print(pHbyDOPlot)
```

### The Effect of pH on DO Concentrations



Insert the following line of code into your R chunk. This will eliminate duplicate measurements on single dates for each site.

```
HawaiiWaterCleanOahu$Station.No<-as.integer(HawaiiWaterCleanOahu$Station.No) ###Removing duplicate measurements
HawaiiWaterCleanOahu2 = HawaiiWaterCleanOahu[order(HawaiiWaterCleanOahu[, 'Date']),]
HawaiiWaterCleanOahu2= HawaiiWaterCleanOahu[!duplicated(HawaiiWaterCleanOahu>Date),]
```

## Mixed Effects Model with Cleaned Data

```
OahuMixed<- lme(data = HawaiiWaterCleanOahu2,
                  DO~Date, #fixed effects model with an interaction term
                  random = ~1|Location) #specifying a random effect

summary(OahuMixed)

## Linear mixed-effects model fit by REML
##  Data: HawaiiWaterCleanOahu2
##      AIC      BIC    logLik
##  723.3878 739.3033 -357.6939
## 
## Random effects:
##   Formula: ~1 | Location
##             (Intercept) Residual
```

```

## StdDev: 0.5260157 0.5518573
##
## Fixed effects: DO ~ Date
##             Value Std.Error DF   t-value p-value
## (Intercept) 8.930163 1.566282 369 5.701505 0.0000
## Date       -0.000230 0.000121 369 -1.899397 0.0583
## Correlation:
##      (Intr)
## Date -0.997
##
## Standardized Within-Group Residuals:
##      Min        Q1        Med        Q3        Max
## -3.838717881 -0.507522085 -0.003601932  0.610439338  2.802921881
##
## Number of Observations: 397
## Number of Groups: 27

```

## ACF

```
#ACF(OahuMixed)
```

0.12-12% of variability associated with time is autocorrelated from previous dates

## Repeated Measures ANOVA Model

```
OahuMixedMod<- lme(data = HawaiiWaterCleanOahu2,
                     DO~Date, #fixed effects
                     random = ~1|Location, #random effect
                     correlation = corAR1(form = ~ Date|Location, value = 0.12),method = "REML")
```

### Summary of OahuMixedMod

```
summary(OahuMixedMod)

## Linear mixed-effects model fit by REML
## Data: HawaiiWaterCleanOahu2
##      AIC      BIC      logLik
## 709.8602 729.7547 -349.9301
##
## Random effects:
## Formula: ~1 | Location
##          (Intercept) Residual
## StdDev: 0.5260025 0.5579336
##
## Correlation Structure: ARMA(1,0)
## Formula: ~Date | Location
## Parameter estimate(s):
##     Phi1
## 0.701675
## Fixed effects: DO ~ Date
```

```

##           Value Std.Error DF   t-value p-value
## (Intercept) 8.468782 1.6999775 369  4.981702 0.0000
## Date        -0.000194 0.0001315 369 -1.477164 0.1405
## Correlation:
##      (Intr)
## Date -0.997
##
## Standardized Within-Group Residuals:
##      Min       Q1       Med       Q3       Max
## -3.78602900 -0.49631030  0.01219029  0.60118603  2.74826863
##
## Number of Observations: 397
## Number of Groups: 27

```

According to the summary of our mixed effects model, the coefficient for the parameter of Date is -0.000194 ( $p=0.14$ ,  $t=-1.47$ ,  $df=369$ ). However, according to the summary, Date is not a significant predictor of PM2.5 concentration because the  $p$ -value for Date is 0.14, which is above 0.05. Thus, there is not a significant trend in DO Concentrations

## Run a Fixed Effects Model with Date as the only predictor

```

OahuFixedMod<- gls(data =HawaiiWaterCleanOahu2,
                     DO~ Date, method="REML")
summary(OahuFixedMod)

## Generalized least squares fit by REML
##   Model: DO ~ Date
##   Data: HawaiiWaterCleanOahu2
##          AIC      BIC    logLik
##  824.6624 836.5991 -409.3312
##
## Coefficients:
##           Value Std.Error t-value p-value
## (Intercept) 9.035947 1.2967383 6.968211 0.0000
## Date        -0.000261 0.0001001 -2.607028 0.0095
##
## Correlation:
##      (Intr)
## Date -1
##
## Standardized residuals:
##      Min       Q1       Med       Q3       Max
## -2.89276470 -0.58541595  0.03111128  0.59653678  3.77943705
##
## Residual standard error: 0.6619924
## Degrees of freedom: 397 total; 395 residual

```

## Compare Mixed Effects Mod to Fixed Effects Mod

```
anova(OahuMixedMod, OahuFixedMod)

##          Model df      AIC      BIC    logLik   Test L.Ratio p-value
## OahuMixedMod     1  5 709.8602 729.7547 -349.9301
## OahuFixedMod     2  3 824.6624 836.5991 -409.3312 1 vs 2 118.8022 <.0001
```

According to the ANOVA test, there is more variability in model structure (error) accounted for by the mixed effects model that includes Location as a random effect. We know this because the AIC score of the OahuMixedMod is 709.86, compared to the OahuFixedMod's AIC score of 824.66. The p-value of <0.0001 indicates that the model fit is significantly different between the two models. Thus, the Mixed Effects model is the best model.

## Add More Parameters, keeping Location as the Random Effect

```
library(lme4)
OahuMixed2<- lme(data = HawaiiWaterCleanOahu2,
                  DO~ Date*Enterococcus*Temperature*Salinity*Turbidity, ####won't let me use pH or time
                  random = ~1|Location) #####R won't let me use week, month, or year as a random
```

## Determine Temporal Autocorrelation in Model

```
#ACF(OahuMixed2)
```

MixedMod-Doesn't WORK, NONE OF THE PARAMETERS ARE SIGNIFICANT

```
#OahuMixedMod2<- lme(data = HawaiiWaterCleanOahu2, DO~ Date*Enterococcus*Temperature*Salinity*Turbidity
#random = ~1/Location,
#correlation = corAR1(form = ~ Date/Location, value = 0.114),
#method = "REML")

#summary(OahuMixedMod2)
```

Research question: Is there a trend over time in DO concentrations by region?

Split Dataset by Region (Use full dataset)

```
HawaiiWaterCleanOahuNorth<- filter(HawaiiWaterCleanOahu, Region=="North")
HawaiiWaterCleanOahuSouth<- filter(HawaiiWaterCleanOahu, Region=="South")
HawaiiWaterCleanOahuEast<- filter(HawaiiWaterCleanOahu, Region=="East")
HawaiiWaterCleanOahuWest<- filter(HawaiiWaterCleanOahu, Region=="West")
```

## Run a Mann Kendall Test for North Oahu

```
library(trend)
mk.test(HawaiiWaterCleanOahuNorth$D0)

##
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuNorth$D0
## z = 3.808, n = 1002, p-value = 0.0001401
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##          S      varS      tau
## 4.029000e+04 1.119384e+08 8.057285e-02
```

## Run a Mann Kendall Test for South Oahu

```
library(trend)
mk.test(HawaiiWaterCleanOahuSouth$D0)

##
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuSouth$D0
## z = 7.1428, n = 2634, p-value = 9.147e-13
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##          S      varS      tau
## 3.219490e+05 2.031599e+09 9.305543e-02
```

## Run a Mann Kendall Test for East Oahu

```
library(trend)
mk.test(HawaiiWaterCleanOahuEast$D0)

##
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuEast$D0
## z = 2.4076, n = 999, p-value = 0.01606
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##          S      varS      tau
## 2.535900e+04 1.109376e+08 5.100554e-02
```

## Run a Mann Kendall Test for West Oahu

```
library(trend)
mk.test(HawaiiWaterCleanOahuWest$D0)
```

```

## 
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuWest$DO
## z = -4.0098, n = 969, p-value = 6.078e-05
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S          varS          tau
## -4.034700e+04 1.012427e+08 -8.630965e-02

```

For North Oahu, the z-value is 3.808, so we see a positive trend in DO concentrations over time. The p-value is 0.00014, so we reject the null hypothesis that the data come from a population with independent realizations and are identically distributed. For South Oahu, the z-value is 7.14, so we see a positive trend in DO concentrations over time. The p-value for Paul Lake is listed as 9.14e-13, so we reject the null hypothesis that the data come from a population with independent realizations and are identically distributed. For East Oahu, the z-value is 2.4. For West Oahu, the z-value is -4.009.

## Pettit's Test for North Oahu

```

pettitt.test(HawaiiWaterCleanOahuNorth$DO)

##
##  Pettitt's test for single change-point detection
##
## data: HawaiiWaterCleanOahuNorth$DO
## U* = 59236, p-value = 1.665e-09
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                               347

```

Because the p-value is <0.05, the change point is significant. Given 1st change point for Peter Lake is 347, we scroll to observation 347 in data set, so first change point occurred in 2004-08-16

## Run a separate Mann-Kendall Test for Each Change Point

```

mk.test(HawaiiWaterCleanOahuNorth$DO[1:346])

##
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuNorth$DO[1:346]
## z = -1.3187, n = 346, p-value = 0.1873
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S          varS          tau
## -2.836000e+03 4.621646e+06 -4.766932e-02
mk.test(HawaiiWaterCleanOahuNorth$DO[347:1002])

```

```

## 
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuNorth$DO[347:1002]
## z = -2.6982, n = 656, p-value = 0.006972
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S          varS          tau
## -1.512900e+04 3.143569e+07 -7.062577e-02

```

p-value for [347:1002] is significant, so run a Pettit's Test

Is there a second change point?

```
pettitt.test(HawaiiWaterCleanOahuNorth$DO[347:1002])
```

```

## 
##  Pettitt's test for single change-point detection
##
## data: HawaiiWaterCleanOahuNorth$DO[347:1002]
## U* = 23298, p-value = 1.988e-05
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                               361

```

347 + 360=707, so look at 707th row->Observation occurred on 2005-08-17

Run another Mann-Kendall for the second change point

Now split dataset into three pieces

```

mk.test(HawaiiWaterCleanOahuNorth$DO[347:706])

## 
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuNorth$DO[347:706]
## z = 3.7634, n = 360, p-value = 0.0001676
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S          varS          tau
## 8.587000e+03 5.205018e+06 1.332158e-01

mk.test(HawaiiWaterCleanOahuNorth$DO[707:1002])

## 
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuNorth$DO[707:1002]
## z = -0.53066, n = 296, p-value = 0.5957
## alternative hypothesis: true S is not equal to 0
## sample estimates:

```

```

##           S          varS          tau
## -9.040000e+02  2.895633e+06 -2.078175e-02

```

If z-score is positive, it's a positive trend. If z-score is negative, it is a negative trend

There is a significant trend in DO concentrations over time at North Oahu for rows:347:706 because the p-value is 0.00016, which is below 0.05.

There is also not a significant trend in DO concentrations over time at North Oahu for rows 707:1002 because the p-value is 0.59, which is above 0.05.

## Pettit's Test for South Oahu

```

pettitt.test(HawaiiWaterCleanOahuSouth$DO)

##
##  Pettitt's test for single change-point detection
##
## data: HawaiiWaterCleanOahuSouth$DO
## U* = 381900, p-value < 2.2e-16
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                               2260

```

Change point is significant bc p<0.05. Given 1st change point for South Oahu is 2260, we scroll to observation 2260 in data set, so first change point occurred in 2005-02-28

## Run separate Mann-Kendall Test for each change point in South Oahu

```

mk.test(HawaiiWaterCleanOahuSouth$DO[1:2259])

##
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuSouth$DO[1:2259]
## z = -1.8701, n = 2259, p-value = 0.06146
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S          varS          tau
## -6.695300e+04  1.281678e+09 -2.631344e-02

mk.test(HawaiiWaterCleanOahuSouth$DO[2260:2634])

##
##  Mann-Kendall trend test
##

```

```

## data: HawaiiWaterCleanOahuSouth$DO[2260:2634]
## z = 3.6784, n = 375, p-value = 0.0002347
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S          varS         tau
## 8.922000e+03 5.881904e+06 1.276411e-01

```

Is there a second change point?

```

pettitt.test(HawaiiWaterCleanOahuSouth$DO[2260:2634])

##
##  Pettitt's test for single change-point detection
##
## data: HawaiiWaterCleanOahuSouth$DO[2260:2634]
## U* = 13508, p-value = 2.036e-09
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                               188

```

The p-value is significant, so there is a second change point.  $2260 + 187 = 2447$ . Look at 2447'th row for second change point, it occurred in 2006-10-12

Run another Mann-Kendall for the second change point

Now split dataset into three pieces

```

mk.test(HawaiiWaterCleanOahuSouth$DO[2260:2446])

##
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuSouth$DO[2260:2446]
## z = -3.5527, n = 187, p-value = 0.0003813
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S          varS         tau
## -3.041000e+03 7.322083e+05 -1.753959e-01

mk.test(HawaiiWaterCleanOahuSouth$DO[2447:2634])

##
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuSouth$DO[2447:2634]
## z = -1.5432, n = 188, p-value = 0.1228
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S          varS         tau
## -1.332000e+03 7.438500e+05 -7.613252e-02

```

If z-score is positive, it's a positive trend. If z-score is negative, it is a negative trend

There is a significant trend in DO concentrations over time at South Oahu for rows:2260:2446 because the p-value is 0.00038, which is below 0.05.

There is also not a significant trend in DO concentrations over time for rows 2447:2634 because the p-value is 0.12, which is above 0.05.

## Petitt Test for East Oahu

```
pettitt.test(HawaiiWaterCleanOahuEast$DO)
```

```
##  
##  Pettitt's test for single change-point detection  
##  
## data: HawaiiWaterCleanOahuEast$DO  
## U* = 38689, p-value = 0.0002471  
## alternative hypothesis: two.sided  
## sample estimates:  
## probable change point at time K  
##                                675
```

Because the p-value is <0.05, the change point is significant. Given 1st change point for East Oahu is 675, we scroll to observation 675 in data set, so first change point occurred in 2006-09-25

## Run separate Mann-Kendall Test for each change point

```
mk.test(HawaiiWaterCleanOahuEast$DO[1:674])
```

```
##  
##  Mann-Kendall trend test  
##  
## data: HawaiiWaterCleanOahuEast$DO[1:674]  
## z = -2.7967, n = 674, p-value = 0.005162  
## alternative hypothesis: true S is not equal to 0  
## sample estimates:  
##           S          varS          tau  
## -1.633100e+04 3.409347e+07 -7.219551e-02  
mk.test(HawaiiWaterCleanOahuEast$DO[675:999])
```

```
##  
##  Mann-Kendall trend test  
##  
## data: HawaiiWaterCleanOahuEast$DO[675:999]  
## z = 1.8734, n = 325, p-value = 0.06101  
## alternative hypothesis: true S is not equal to 0  
## sample estimates:  
##           S          varS          tau  
## 3.668000e+03 3.831328e+06 6.985628e-02
```

There are no more change points because p-value is >0.05 for interval [675:999]

## Petitt Test for West Oahu

```
pettitt.test(HawaiiWaterCleanOahuWest$DO)

##
##  Pettitt's test for single change-point detection
##
## data: HawaiiWaterCleanOahuWest$DO
## U* = 53608, p-value = 1.2e-08
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                               663
```

Because the p-value is <0.05, the change point is significant. Given 1st change point for West Oahu is 663, we scroll to observation 663 in data set, so first change point occurred in 2003-12-03

## Run separate Mann-Kendall Test for each change point

```
mk.test(HawaiiWaterCleanOahuWest$DO[1:662])

##
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuWest$DO[1:662]
## z = 0.67807, n = 662, p-value = 0.4977
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S      varS      tau
## 3.855000e+03 3.230499e+07 1.767791e-02

mk.test(HawaiiWaterCleanOahuWest$DO[663:969])

##
##  Mann-Kendall trend test
##
## data: HawaiiWaterCleanOahuWest$DO[663:969]
## z = 5.0599, n = 307, p-value = 4.194e-07
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##           S      varS      tau
## 9.095000e+03 3.230146e+06 1.942265e-01
```

p-value for [663:969] is significant, so there is a second change point

What is the second change point?

```
pettitt.test(HawaiiWaterCleanOahuWest$DO[663:969])  
  
##  
##  Pettitt's test for single change-point detection  
##  
## data: HawaiiWaterCleanOahuWest$DO[663:969]  
## U* = 9574, p-value = 1.183e-08  
## alternative hypothesis: two.sided  
## sample estimates:  
## probable change point at time K  
##                               165  
  
663+164=827, second change point occurred on 2004-10-20
```

Run another Mann-Kendall for the third change point (if it exists)

Now split dataset into three pieces

```
mk.test(HawaiiWaterCleanOahuWest$DO[663:826])  
  
##  
##  Mann-Kendall trend test  
##  
## data: HawaiiWaterCleanOahuWest$DO[663:826]  
## z = -0.08533, n = 164, p-value = 0.932  
## alternative hypothesis: true S is not equal to 0  
## sample estimates:  
##           S          varS          tau  
## -6.100000e+01  4.944290e+05 -4.579434e-03  
  
mk.test(HawaiiWaterCleanOahuWest$DO[827:969])  
  
##  
##  Mann-Kendall trend test  
##  
## data: HawaiiWaterCleanOahuWest$DO[827:969]  
## z = -0.39801, n = 143, p-value = 0.6906  
## alternative hypothesis: true S is not equal to 0  
## sample estimates:  
##           S          varS          tau  
## -2.29000e+02  3.28165e+05 -2.26488e-02
```

There is not a third changepoint

## Time Series Graph of North Oahu

```
NorthOahuTimeSeries<-ggplot(HawaiiWaterCleanOahuNorth, aes(x = Date, y = DO)) +  
  geom_point(alpha=1) +
```

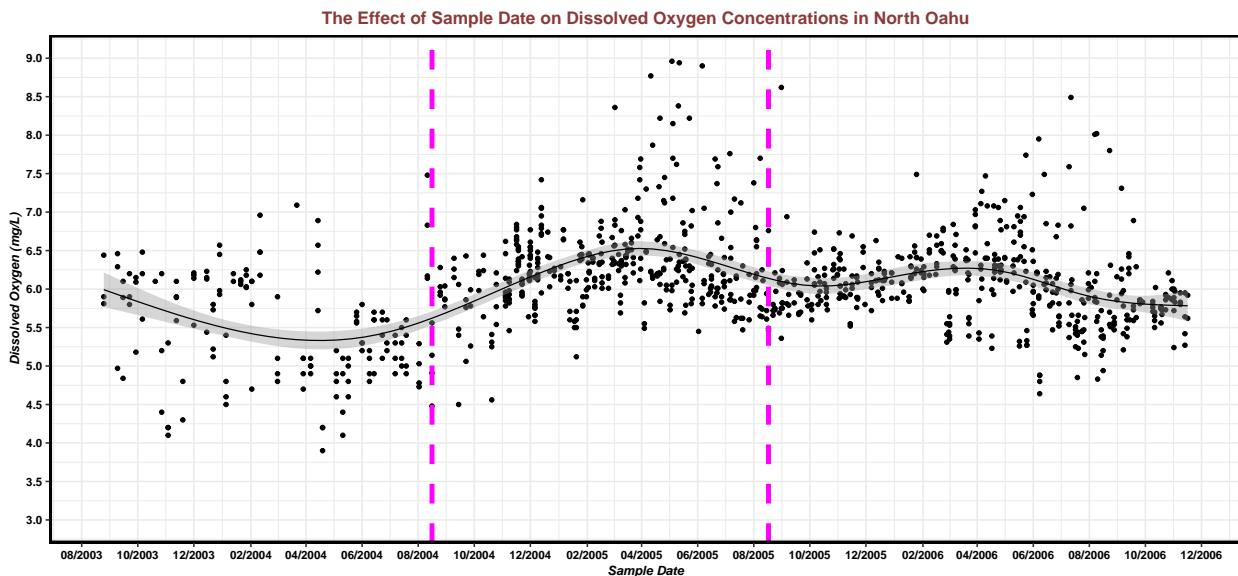
```

geom_vline(xintercept=as.Date("2004-08-16"), color="253494", origin= "1970-01-01", lty=2, lwd=2) + ###First Vline
geom_vline(xintercept=as.Date("2005-08-17"), color="253494", origin= "1970-01-01", lty=2, lwd=2) + ###Second Vline

geom_smooth(aes(x = Date, y = DO, span=0.1), color="black", linetype=1, size=0.5) +
  labs(title="The Effect of Sample Date on Dissolved Oxygen Concentrations in North Oahu",
       x="Sample Date",
       y="Dissolved Oxygen (mg/L)") +
  scale_x_date(labels = date_format("%m/%Y"), breaks = date_breaks("2 month")) +
  scale_y_continuous(limits=c(3,9), breaks=seq(3,9, by = 0.5))

```

NorthOahuTimeSeries



## Effect of Sample Date on Dissolved Oxygen Concentration

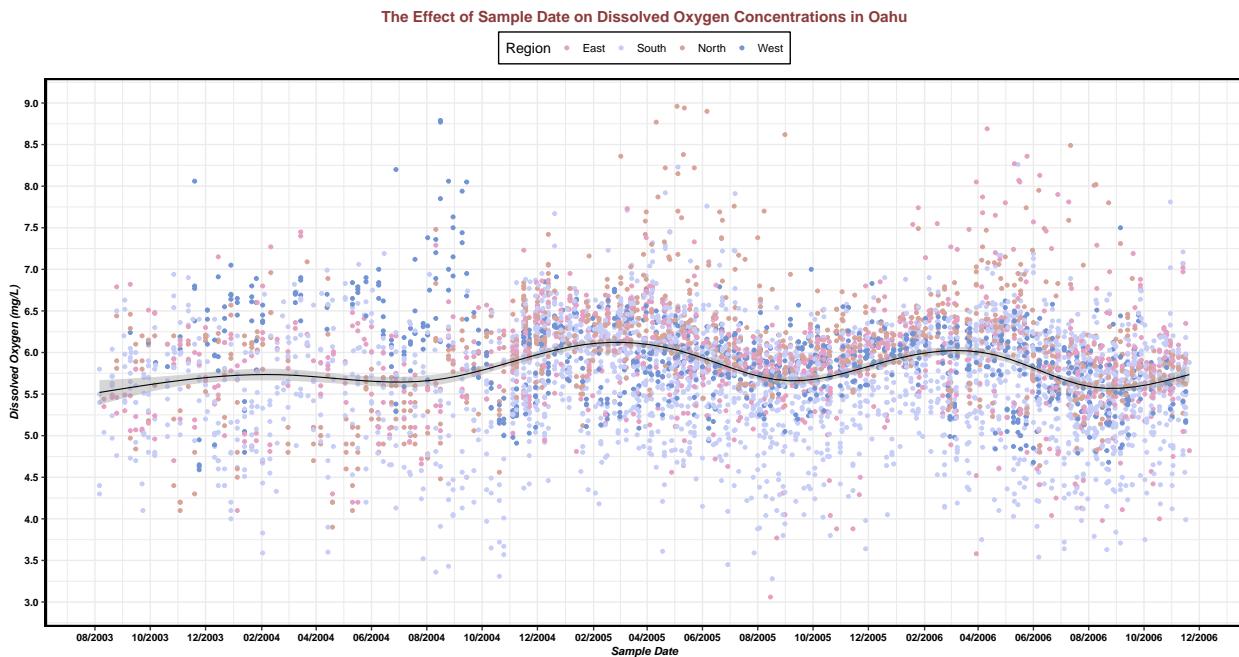
```

library(wesanderson)
OahuDOPPlot<-ggplot(HawaiiWaterCleanOahu, aes(x = Date, y = DO, color = Region)) +
  geom_point(alpha=1) +
  scale_color_manual(values=wes_palette(name="GrandBudapest2")) +
  ##geom_vline(xintercept=as.Date("2004-08-16"), color="253494", origin= "1970-01-01", lty=2) + ###First Vline
  ## geom_vline(xintercept=as.Date("2005-08-17"), color="253494", origin= "1970-01-01", lty=2)  ###Second Vline

  geom_smooth(aes(x = Date, y = DO, span=0.1), color="black", linetype=1, size=0.5) +
  labs(title="The Effect of Sample Date on Dissolved Oxygen Concentrations in Oahu",
       x="Sample Date",
       y="Dissolved Oxygen (mg/L)") +
  scale_x_date(labels = date_format("%m/%Y"), breaks = date_breaks("2 month")) +
  scale_y_continuous(limits=c(3,9), breaks=seq(3,9, by = 0.5))

```

## OahuDOPPlot

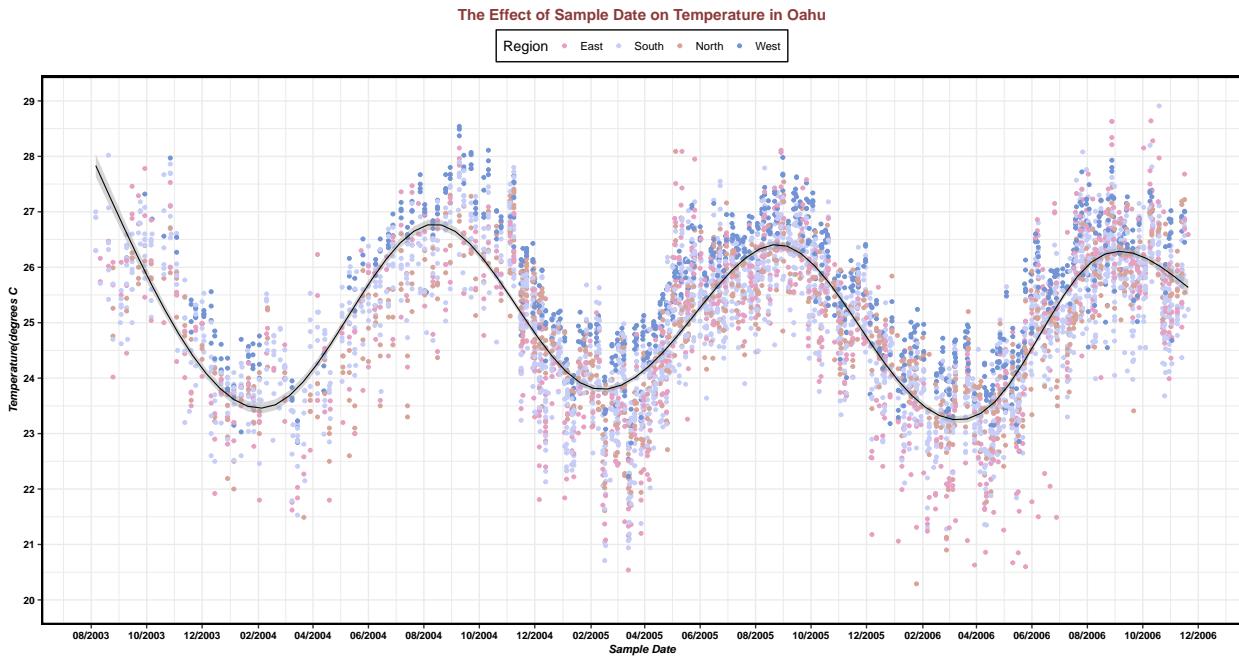


## Effect of Sample Date on Temperature in Oahu

```
library(wesanderson)
OahuTemperaturePlot<-ggplot(HawaiiWaterCleanOahu, aes(x = Date, y = Temperature, color = Region)) +
  geom_point(alpha=1) +
  scale_color_manual(values=wes_palette(name="GrandBudapest2")) +
  ##geom_vline(xintercept=as.Date("2004-08-16"),color="253494", origin= "1970-01-01", lty=2) + ###First
## geom_vline(xintercept=as.Date("2005-08-17"), color="253494", origin= "1970-01-01", lty=2)  ###Second

  geom_smooth(aes(x = Date, y = Temperature, span=0.1), color="black", linetype=1, size=0.5) +
  labs(title="The Effect of Sample Date on Temperature in Oahu",
       x="Sample Date",
       y="Temperature(degrees C)") +
  scale_x_date(labels = date_format("%m/%Y"), breaks = date_breaks("2 month")) +
  scale_y_continuous(limits=c(20,29), breaks=seq(20, 29, by =1))
```

```
OahuTemperaturePlot
```



## Effect of Sample Date on pH in Oahu

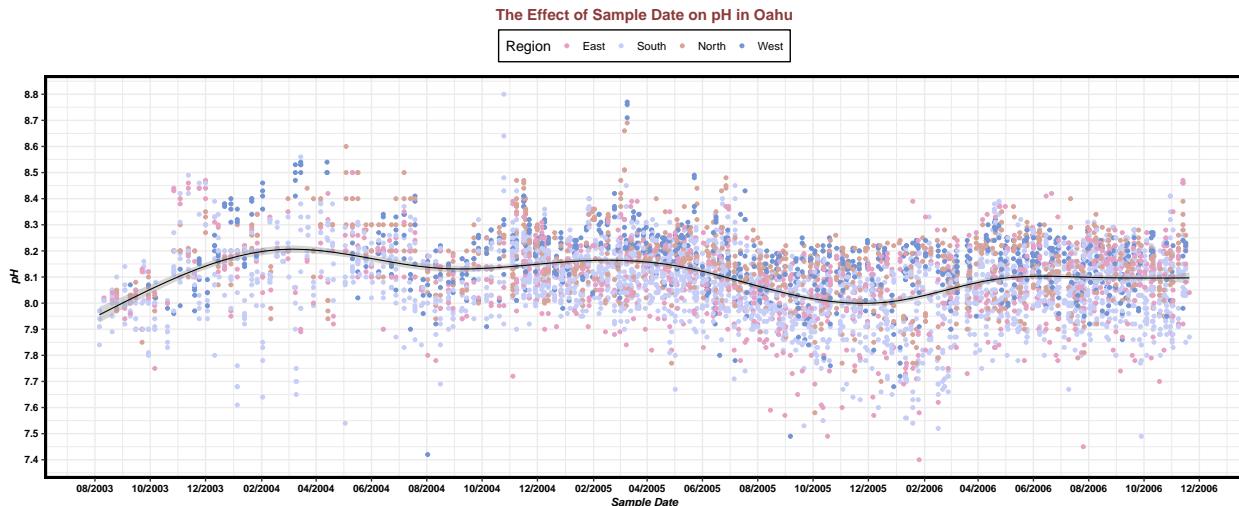
```

library(wesanderson)
OahupHPlot<-ggplot(HawaiiWaterCleanOahu, aes(x = Date, y =pH, color = Region)) +
  geom_point(alpha=1) +
  scale_color_manual(values=wes_palette(name="GrandBudapest2")) +
  ##geom_vline(xintercept=as.Date("2004-08-16"),color="253494", origin= "1970-01-01", lty=2) + ###First
## geom_vline(xintercept=as.Date("2005-08-17"), color="253494", origin= "1970-01-01", lty=2)  ###Second

  geom_smooth(aes(x = Date, y =pH, span=0.1), color="black", linetype=1, size=0.5) +
  labs(title="The Effect of Sample Date on pH in Oahu",
       x="Sample Date",
       y="pH") +
  scale_x_date(labels = date_format("%m/%Y"), breaks = date_breaks("2 month")) +
  scale_y_continuous(limits=c(7.4,8.8), breaks=seq(7.4, 8.8, by = 0.1))

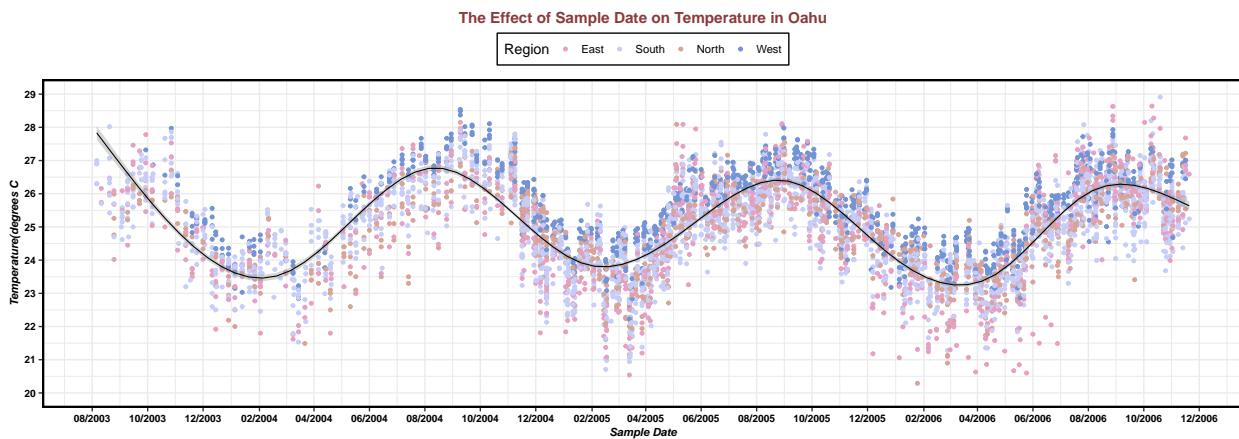
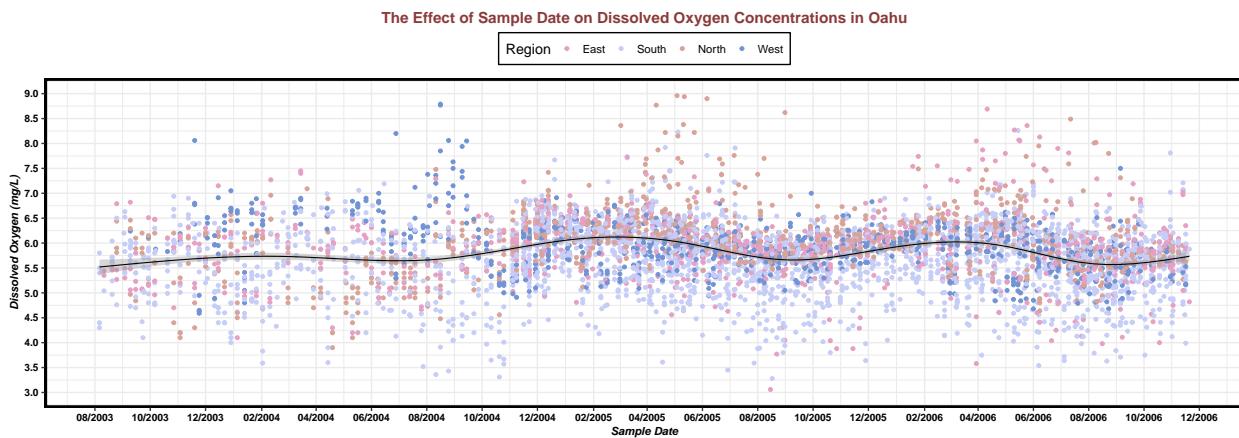
OahupHPlot

```



## Comparison of DO and Temperature over Time in Oahu

```
grid.arrange(OahuDOPlot, OahuTemperaturePlot, nrow=2)
```



Cold water can hold more dissolved oxygen than warm water. In winter and early spring, when the water temperature is low, the dissolved oxygen concentration is high. In summer and fall, when the water temperature is high, the dissolved-oxygen concentration is low.

5. Is there equal variance among the publication years for each chemical? Hint: var.test is not the correct function.

## Summary and Conclusions