# Linear Regression, GAM Models, and RMSE Comparison

2025-05-05

## Set Up Data

```r
# Get data
wineData <- read.csv("~/Downloads/wine-quality-white-and-red.csv")

wineData$type <- as.factor(wineData$type)

# Split to training and testing data
set.seed(1)
train = sample(1:nrow(wineData), 0.7 * nrow(wineData))
train_data = wineData[train, ]
test_data = wineData[-train, ]
```

## Best Subset Selection

```r
regfit.best=regsubsets(quality~.,data=train_data, nvmax=12)
test.mat=model.matrix(quality~.,data=test_data) # create an X matrix of test data
val.errors=rep(NA,19)
for(i in 1:12){
   coefi=coef(regfit.best,id=i)
   pred=test.mat[,names(coefi)]%*%coefi
   val.errors[i] <- mean((test_data$quality - pred)^2)
}

summary(regfit.best)
```

```
## Subset selection object
## Call: regsubsets.formula(quality ~ ., data = train_data, nvmax = 12)
## 12 Variables  (and intercept)
##                     Forced in Forced out
## typewhite              FALSE      FALSE
## fixed.acidity          FALSE      FALSE
## volatile.acidity       FALSE      FALSE
## citric.acid            FALSE      FALSE
## residual.sugar         FALSE      FALSE
## chlorides              FALSE      FALSE
## free.sulfur.dioxide    FALSE      FALSE
## total.sulfur.dioxide   FALSE      FALSE
## density                FALSE      FALSE
## pH                     FALSE      FALSE
```

```
## sulphates                        FALSE      FALSE
## alcohol                          FALSE      FALSE
## 1 subsets of each size up to 12
## Selection Algorithm: exhaustive
##          typewhite fixed.acidity volatile.acidity citric.acid residual.sugar
## 1  ( 1 )  " "       " "           " "              " "         " "
## 2  ( 1 )  " "       " "           "*"              " "         " "
## 3  ( 1 )  " "       " "           "*"              " "         " "
## 4  ( 1 )  " "       " "           "*"              " "         "*"
## 5  ( 1 )  " "       " "           "*"              " "         "*"
## 6  ( 1 )  " "       " "           "*"              " "         "*"
## 7  ( 1 )  " "       " "           "*"              " "         "*"
## 8  ( 1 )  "*"       " "           "*"              " "         "*"
## 9  ( 1 )  "*"       "*"           "*"              " "         "*"
## 10 ( 1 )  "*"       "*"           "*"              " "         "*"
## 11 ( 1 )  "*"       "*"           "*"              " "         "*"
## 12 ( 1 )  "*"       "*"           "*"              "*"         "*"
##          chlorides free.sulfur.dioxide total.sulfur.dioxide density pH
## 1  ( 1 )  " "       " "                 " "                  " "     " "
## 2  ( 1 )  " "       " "                 " "                  " "     " "
## 3  ( 1 )  " "       " "                 " "                  " "     " "
## 4  ( 1 )  " "       " "                 " "                  " "     " "
## 5  ( 1 )  " "       "*"                 "*"                  " "     " "
## 6  ( 1 )  " "       "*"                 "*"                  " "     " "
## 7  ( 1 )  "*"       "*"                 "*"                  " "     " "
## 8  ( 1 )  " "       "*"                 "*"                  "*"     " "
## 9  ( 1 )  " "       "*"                 " "                  "*"     "*"
## 10 ( 1 )  " "       "*"                 "*"                  "*"     "*"
## 11 ( 1 )  "*"       "*"                 "*"                  "*"     "*"
## 12 ( 1 )  "*"       "*"                 "*"                  "*"     "*"
##          sulphates alcohol
## 1  ( 1 )  " "       "*"
## 2  ( 1 )  " "       "*"
## 3  ( 1 )  "*"       "*"
## 4  ( 1 )  "*"       "*"
## 5  ( 1 )  " "       "*"
## 6  ( 1 )  "*"       "*"
## 7  ( 1 )  "*"       "*"
## 8  ( 1 )  "*"       "*"
## 9  ( 1 )  "*"       "*"
## 10 ( 1 )  "*"       "*"
## 11 ( 1 )  "*"       "*"
## 12 ( 1 )  "*"       "*"
```

```r
which.min(val.errors)
```

```
## [1] 12
```

```r
coef(regfit.best,12)
```

```
##         (Intercept)            typewhite          fixed.acidity
##        1.190237e+02         -3.734753e-01           9.607001e-02
##      volatile.acidity          citric.acid         residual.sugar
```

```
##      -1.381747e+00         -1.535991e-02        6.575127e-02
##        chlorides  free.sulfur.dioxide total.sulfur.dioxide
##      -8.970624e-01         5.661084e-03       -1.625646e-03
##          density                   pH            sulphates
##      -1.184895e+02         5.766522e-01        6.741088e-01
##          alcohol
##       2.107583e-01
```
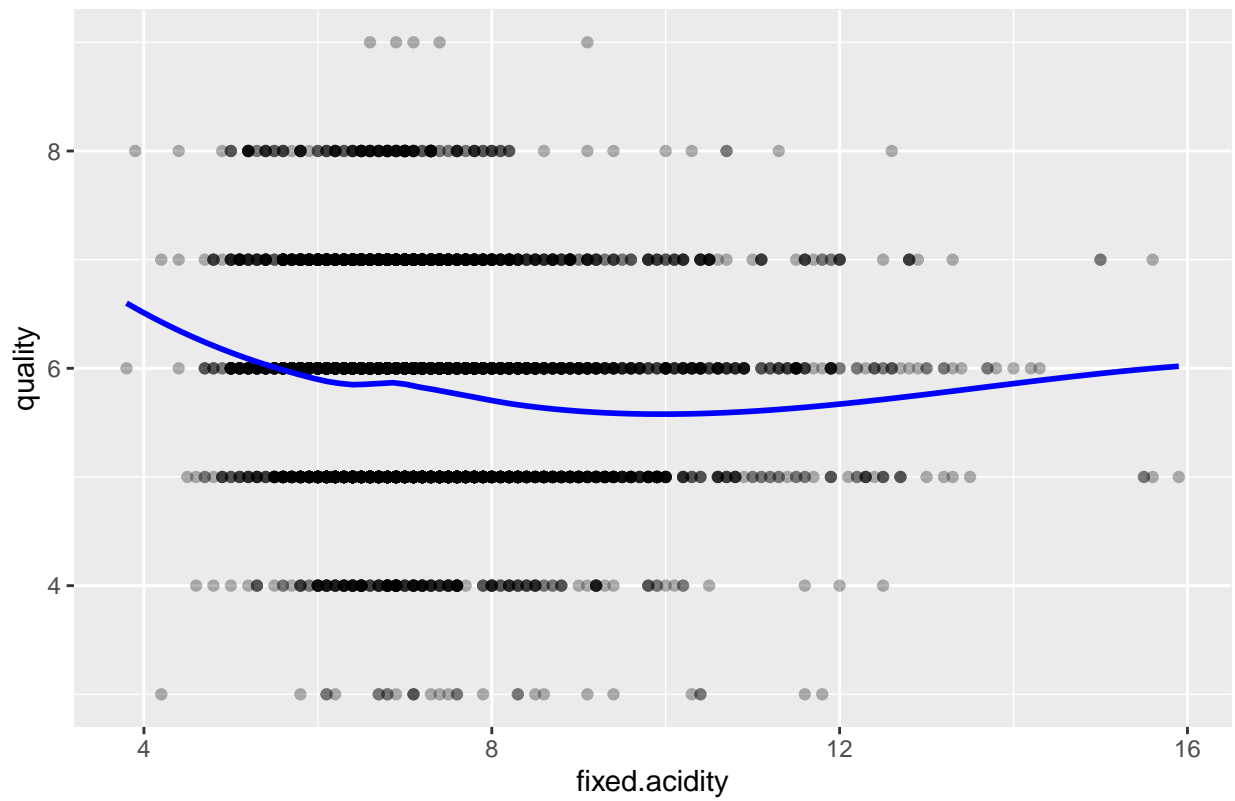
## Explore Predictors for GAM

```r
library(ggplot2)

# Continuous variables
predictors <- c("fixed.acidity", "volatile.acidity", "citric.acid", "residual.sugar",
                "chlorides", "free.sulfur.dioxide", "total.sulfur.dioxide",
                "density", "pH", "sulphates", "alcohol")

# Plot each predictor against quality
for (var in predictors) {
  p <- ggplot(wineData, aes_string(x = var, y = "quality")) +
    geom_point(alpha = 0.3) +
    geom_smooth(method = "loess", se = FALSE, color = "blue") +
    labs(title = paste("Quality vs", var))
  print(p)
}
```

```
## Warning: 'aes_string()' was deprecated in ggplot2 3.0.0.
## i Please use tidy evaluation idioms with 'aes()'.
## i See also 'vignette("ggplot2-in-packages")' for more information.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```
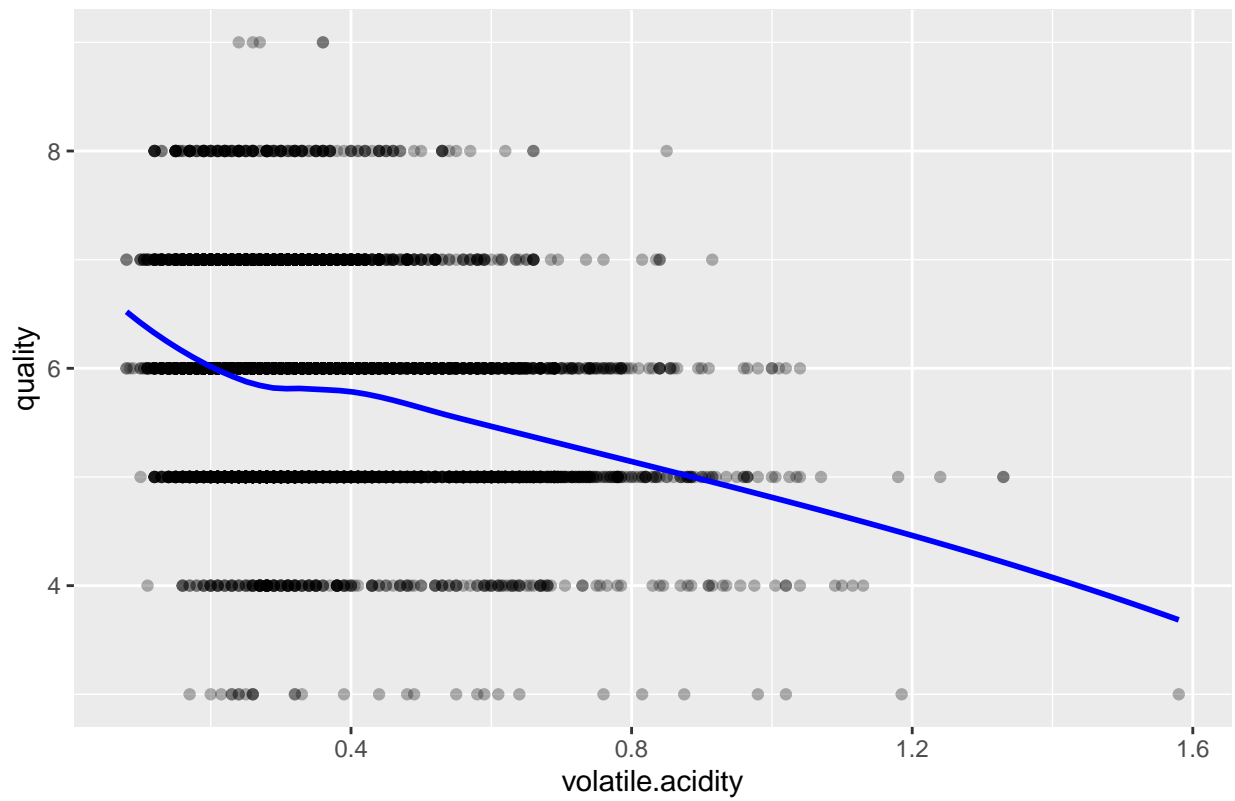
```
## 'geom_smooth()' using formula = 'y ~ x'
```
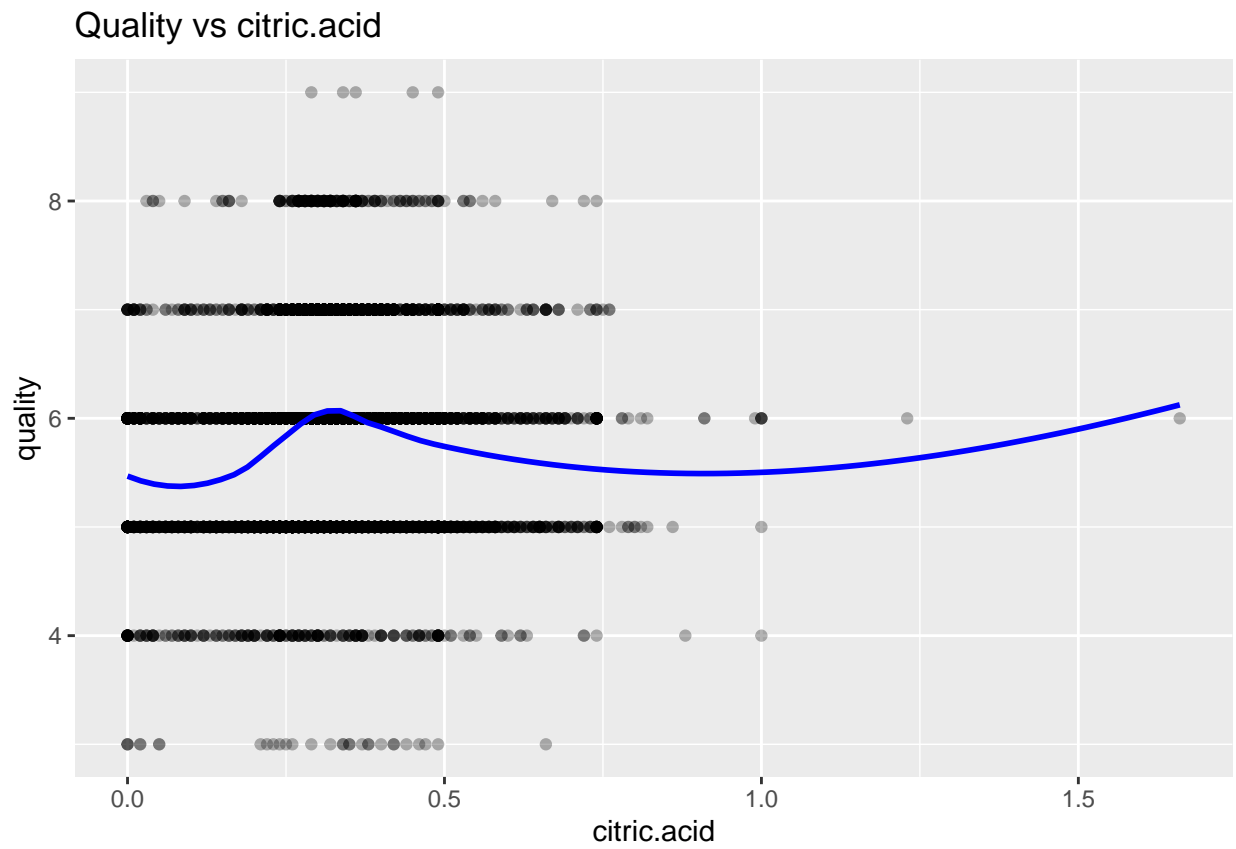
Quality vs fixed.acidity

```
## 'geom_smooth()' using formula = 'y ~ x'
```
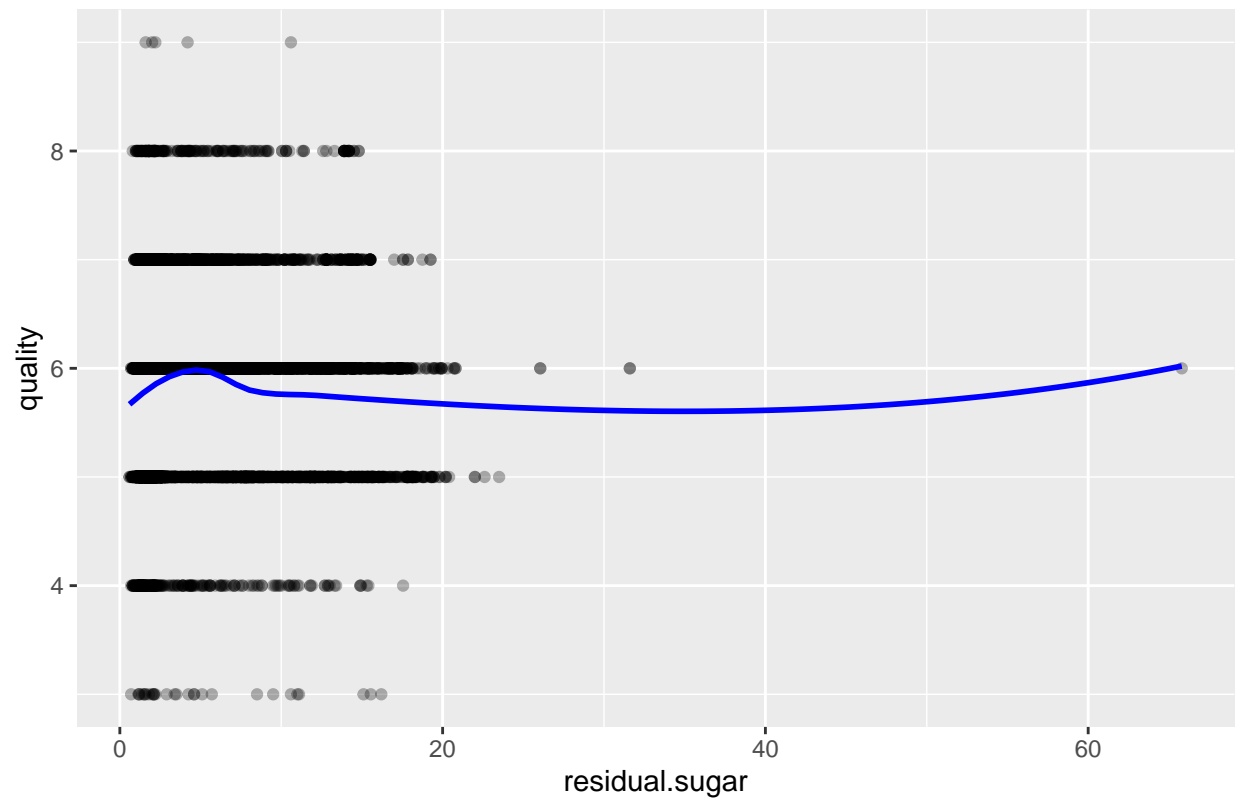
## Quality vs volatile.acidity



```
## `geom_smooth()` using formula = 'y ~ x'
```
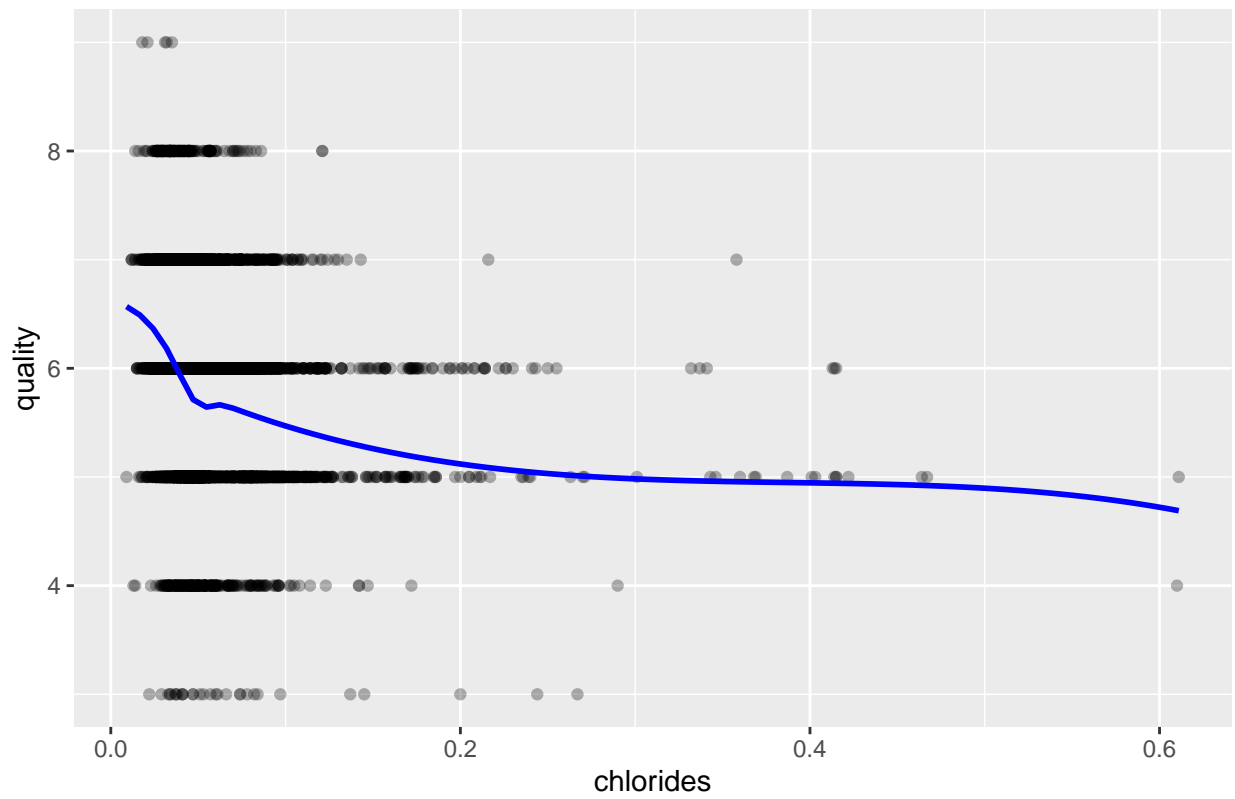
Quality vs citric.acid

```
## 'geom_smooth()' using formula = 'y ~ x'
```
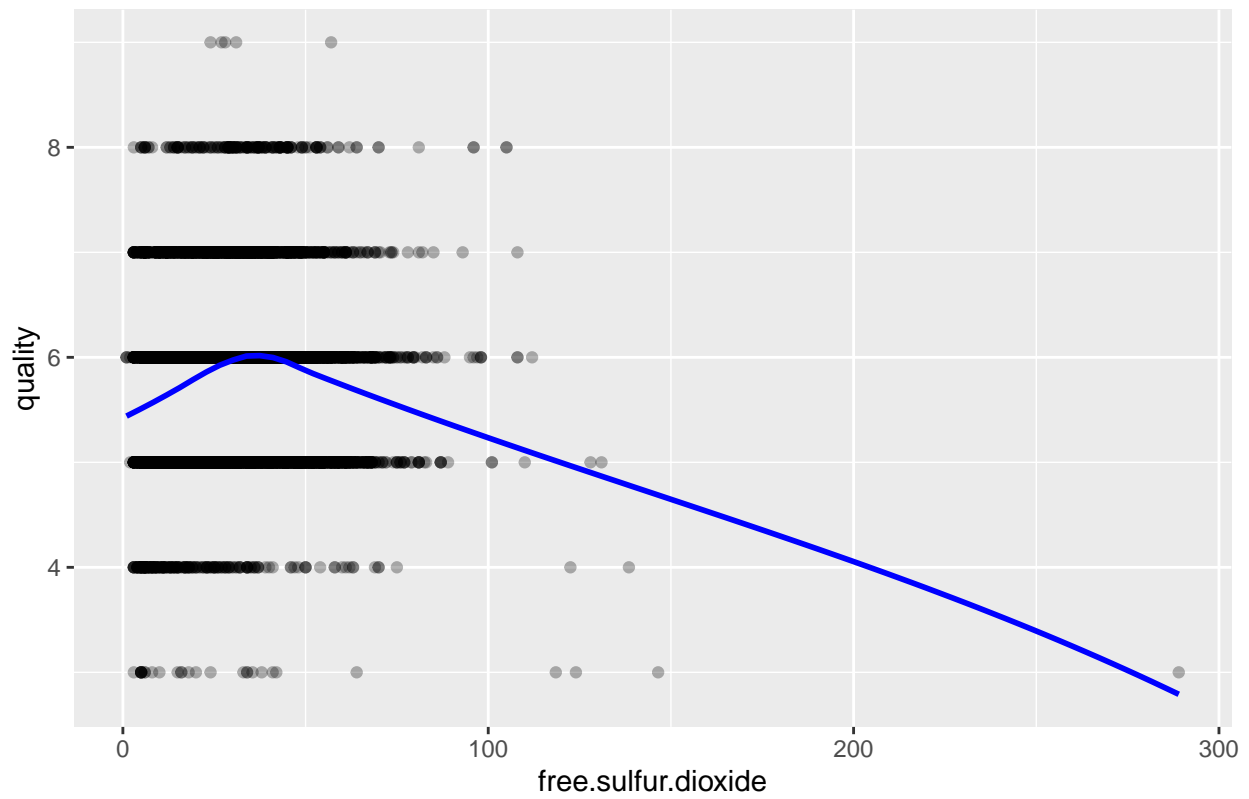
## Quality vs residual.sugar



```
## 'geom_smooth()' using formula = 'y ~ x'
```
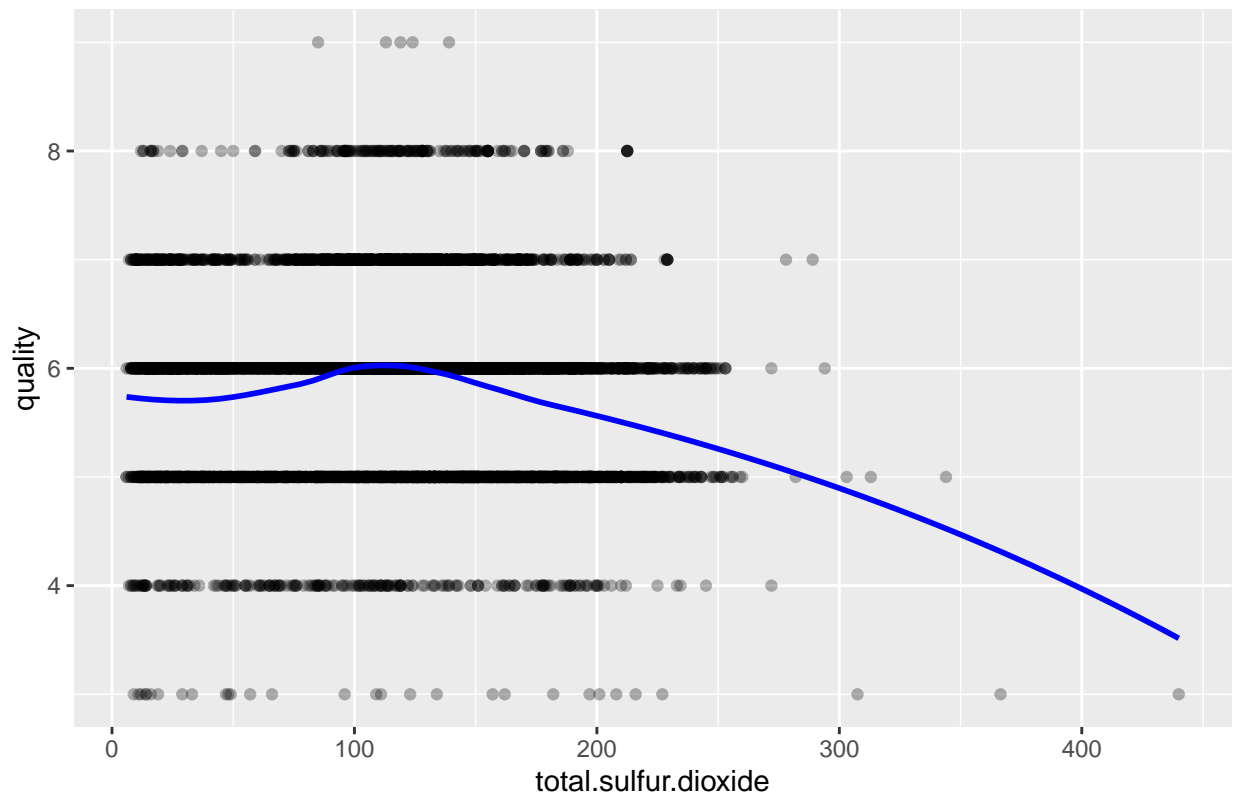
Quality vs chlorides

```
## 'geom_smooth()' using formula = 'y ~ x'
```

## Quality vs free.sulfur.dioxide



```
## 'geom_smooth()' using formula = 'y ~ x'
```

## Quality vs total.sulfur.dioxide



```
## 'geom_smooth()' using formula = 'y ~ x'
```

# Quality vs density



```
## `geom_smooth()` using formula = 'y ~ x'
```
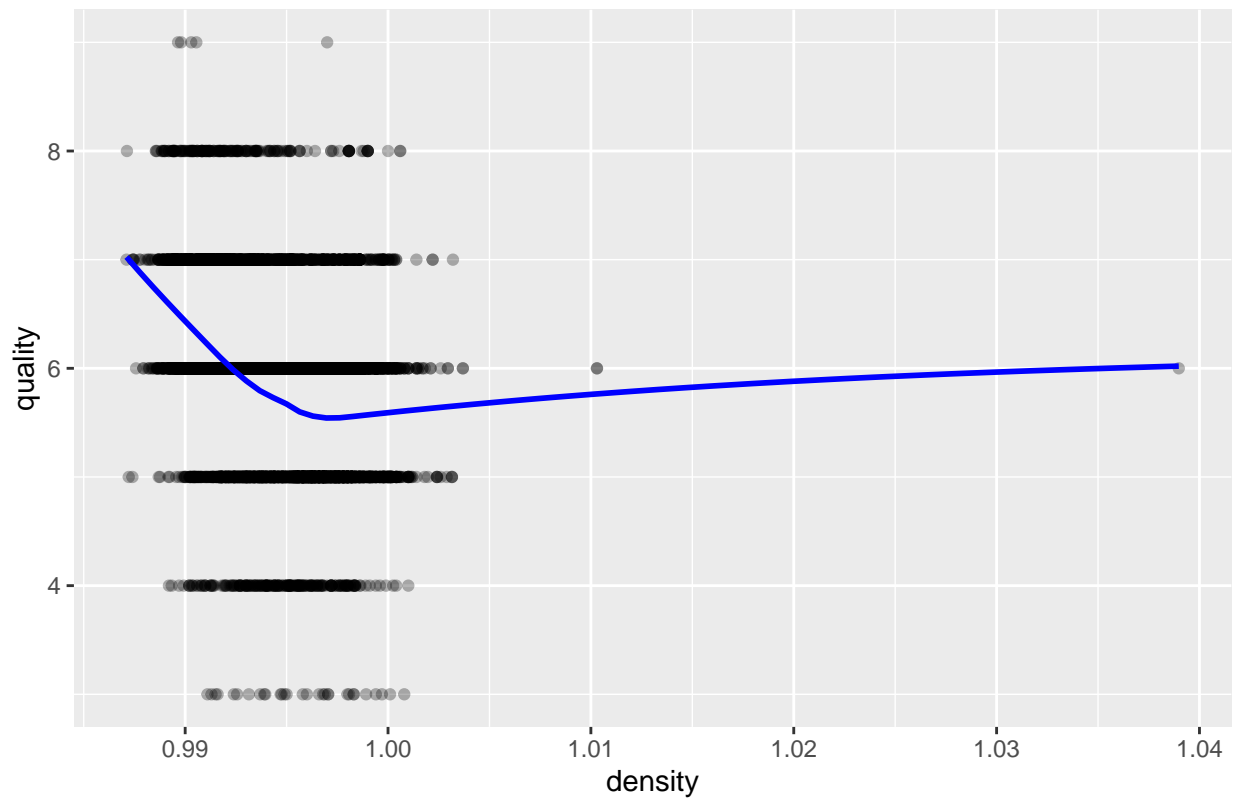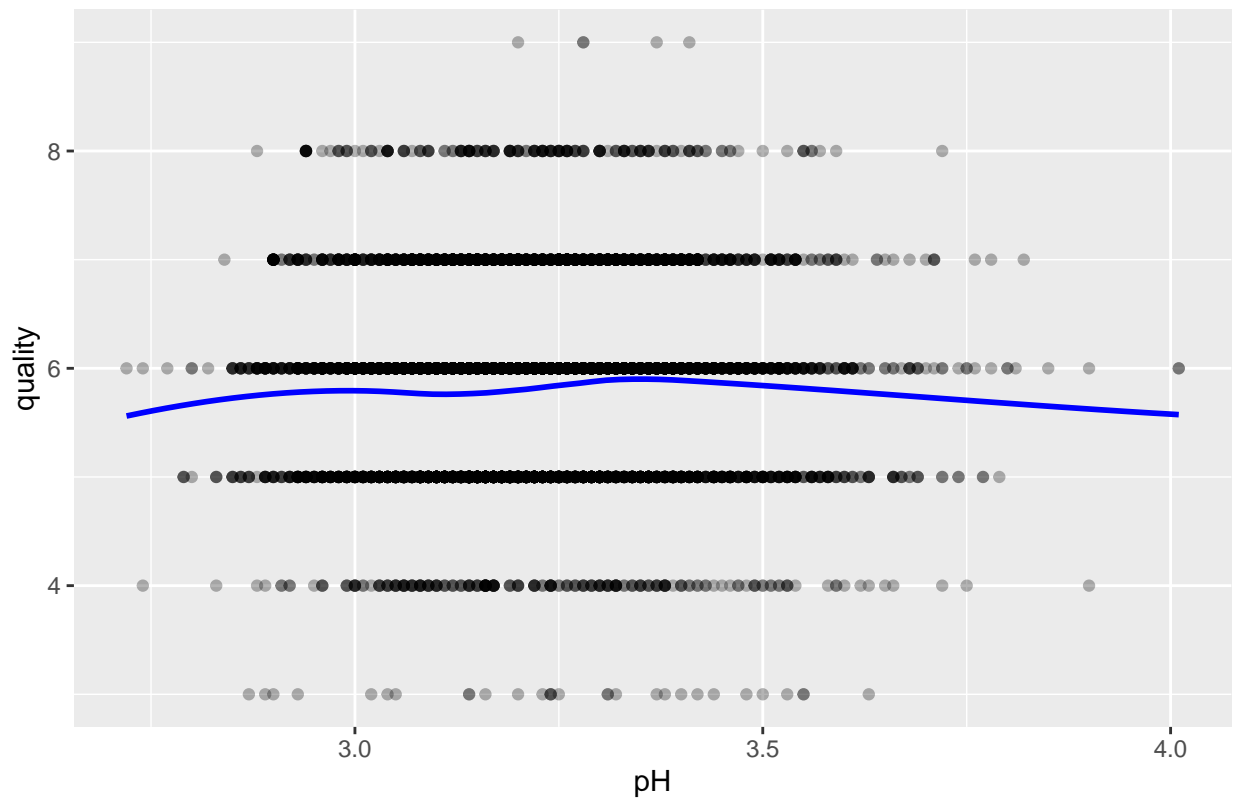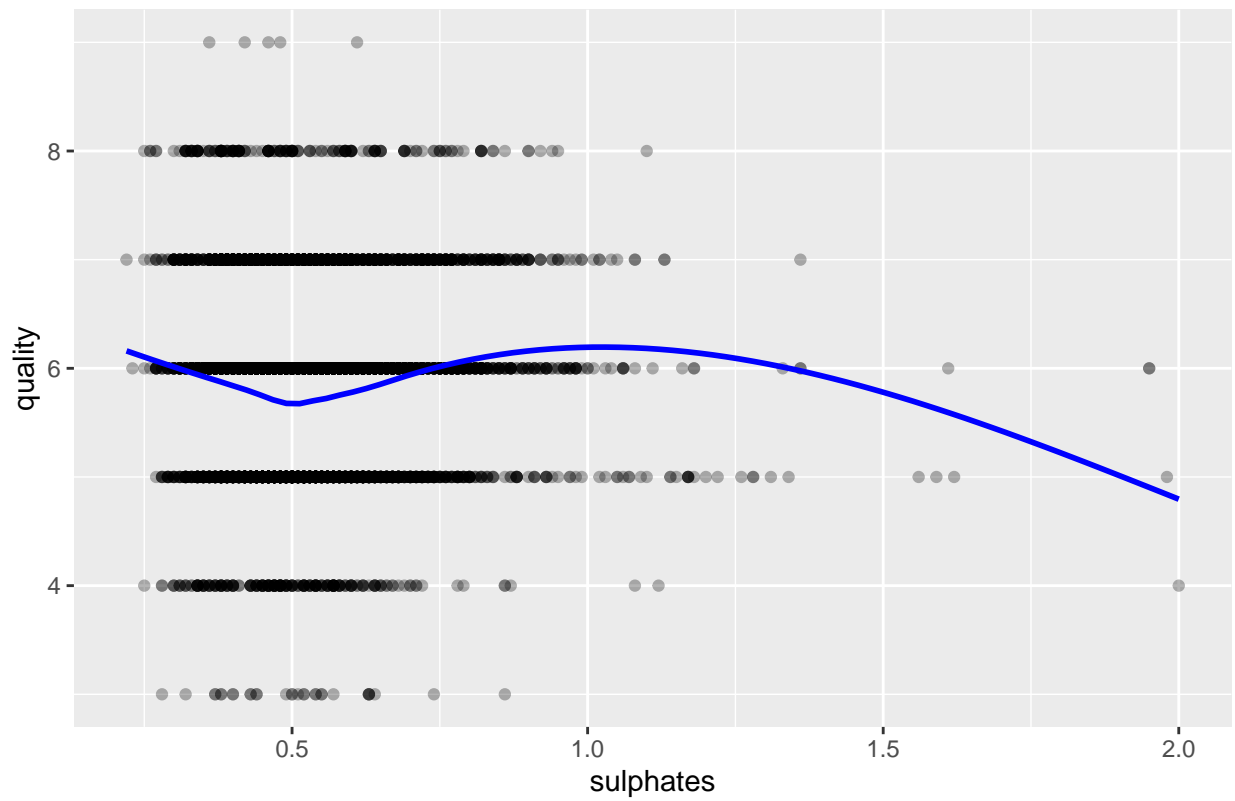
## `geom_smooth()` using formula = 'y ~ x'

Quality vs sulphates

```
## 'geom_smooth()' using formula = 'y ~ x'
```

## Quality vs alcohol



```
# Categorical: type
ggplot(wineData, aes(x = type, y = quality)) +
  geom_boxplot() +
  labs(title = "Quality vs Type")
```

## Quality vs Type



## GAM Model 1 Using Train/Test Split

```r
library(gam)
```

```
## Loading required package: splines
```

```
## Loading required package: foreach
```

```
## Loaded gam 1.22-5
```

```
##
## Attaching package: 'gam'
```

```
## The following objects are masked from 'package:mgcv':
##
##     gam, gam.control, gam.fit, s
```

```r
# Example GAM model using smoothing for a few predictors
gam_model1 <- gam(quality ~
                    s(fixed.acidity) +
                    s(volatile.acidity) +
                    s(citric.acid) +
```

```
                    s(chlorides) +
                     s(residual.sugar) +
                    s(free.sulfur.dioxide) +
                     s(total.sulfur.dioxide) +
                     s(density) +
                     s(sulphates) +
                     s(alcohol) + s(pH), data = train_data)

summary(gam_model1)
```

```
##
## Call: gam(formula = quality ~ s(fixed.acidity) + s(volatile.acidity) +
##      s(citric.acid) + s(chlorides) + s(residual.sugar) + s(free.sulfur.dioxide) +
##      s(total.sulfur.dioxide) + s(density) + s(sulphates) + s(alcohol) +
##      s(pH), data = train_data)
## Deviance Residuals:
##     Min       1Q  Median       3Q      Max
## -3.5824 -0.4591 -0.0195  0.4492   2.4936
##
## (Dispersion Parameter for gaussian family taken to be 0.5009)
##
##     Null Deviance: 3461.128 on 4546 degrees of freedom
## Residual Deviance: 2254.917 on 4502 degrees of freedom
## AIC: 9806.769
##
## Number of Local Scoring Iterations: NA
##
## Anova for Parametric Effects
##                          Df  Sum Sq Mean Sq  F value     Pr(>F)
## s(fixed.acidity)          1   14.10  14.097  28.1449 1.180e-07 ***
## s(volatile.acidity)       1  156.02 156.022 311.5014 < 2.2e-16 ***
## s(citric.acid)            1    0.01   0.011   0.0223    0.8814
## s(chlorides)              1   34.46  34.465  68.8101 < 2.2e-16 ***
## s(residual.sugar)         1   64.10  64.102 127.9807 < 2.2e-16 ***
## s(free.sulfur.dioxide)    1    0.03   0.030   0.0597    0.8070
## s(total.sulfur.dioxide)   1  165.02 165.024 329.4739 < 2.2e-16 ***
## s(density)                1  229.84 229.839 458.8803 < 2.2e-16 ***
## s(sulphates)              1  109.76 109.758 219.1344 < 2.2e-16 ***
## s(alcohol)                1  150.95 150.947 301.3690 < 2.2e-16 ***
## s(pH)                     1    8.22   8.223  16.4167 5.169e-05 ***
## Residuals              4502 2254.92   0.501
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Anova for Nonparametric Effects
##                       Npar Df Npar F      Pr(F)
## (Intercept)
## s(fixed.acidity)            3  3.705   0.011185 *
## s(volatile.acidity)         3 12.581  3.450e-08 ***
## s(citric.acid)              3  7.803  3.414e-05 ***
## s(chlorides)                3  2.763   0.040500 *
## s(residual.sugar)           3 12.193  6.056e-08 ***
## s(free.sulfur.dioxide)      3 53.123  < 2.2e-16 ***
```

16

```
## s(total.sulfur.dioxide)      3 20.407 3.939e-13 ***
## s(density)                    3  4.033  0.007099 **
## s(sulphates)                  3  2.151  0.091701 .
## s(alcohol)                    3 13.383 1.077e-08 ***
## s(pH)                         3  8.208 1.909e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```r
# Predict
gam_pred <- predict(gam_model1, newdata = test_data)

# RMSE
mse <- mean((gam_pred - test_data$quality)^2)
rmse <- sqrt(mse)
cat("RMSE of GAM1: ", rmse)
```

```
## RMSE of GAM1:  0.7220457
```

## GAM Model 2 Using Train/Test Split

```r
library(gam)

# Example GAM model using smoothing for a few predictors
gam_model2 <- gam(quality ~
                        fixed.acidity +
                    volatile.acidity +
                    s(citric.acid) +
                    residual.sugar +
                    s(chlorides) +
                    free.sulfur.dioxide +
                    total.sulfur.dioxide +
                    density +
                    s(sulphates) +
                    s(alcohol) + s(pH),
                     data = train_data)

summary(gam_model2)
```

```
##
## Call: gam(formula = quality ~ fixed.acidity + volatile.acidity + s(citric.acid) +
##     residual.sugar + s(chlorides) + free.sulfur.dioxide + total.sulfur.dioxide +
##     density + s(sulphates) + s(alcohol) + s(pH), data = train_data)
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3.60483 -0.47131 -0.02554  0.44413  2.59545
##
## (Dispersion Parameter for gaussian family taken to be 0.5276)
##
##      Null Deviance: 3461.128 on 4546 degrees of freedom
## Residual Deviance: 2384.858 on 4520 degrees of freedom
## AIC: 10025.52
```

```
##
## Number of Local Scoring Iterations: NA
##
## Anova for Parametric Effects
##                      Df  Sum Sq Mean Sq  F value    Pr(>F)
## fixed.acidity          1   19.73  19.727  37.3892 1.049e-09 ***
## volatile.acidity       1  190.77 190.767 361.5598 < 2.2e-16 ***
## s(citric.acid)         1    0.46   0.457   0.8669    0.3519
## residual.sugar         1   46.43  46.428  87.9941 < 2.2e-16 ***
## s(chlorides)           1   42.38  42.383  80.3286 < 2.2e-16 ***
## free.sulfur.dioxide    1    0.19   0.192   0.3641    0.5462
## total.sulfur.dioxide   1  122.54 122.538 232.2448 < 2.2e-16 ***
## density                1  236.93 236.926 449.0442 < 2.2e-16 ***
## s(sulphates)           1  115.27 115.271 218.4721 < 2.2e-16 ***
## s(alcohol)             1  175.43 175.427 332.4843 < 2.2e-16 ***
## s(pH)                  1   12.39  12.393  23.4886 1.299e-06 ***
## Residuals           4520 2384.86   0.528
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Anova for Nonparametric Effects
##                     Npar Df  Npar F     Pr(F)
## (Intercept)
## fixed.acidity
## volatile.acidity
## s(citric.acid)           3   6.3365 0.0002771 ***
## residual.sugar
## s(chlorides)             3   2.6740 0.0456893 *
## free.sulfur.dioxide
## total.sulfur.dioxide
## density
## s(sulphates)             3   2.0634 0.1028790
## s(alcohol)               3  13.6932 6.868e-09 ***
## s(pH)                    3   9.2017 4.567e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```r
# Predict
gam_pred2 <- predict(gam_model2, newdata = test_data)

# RMSE
mse <- mean((gam_pred2 - test_data$quality)^2)
rmse <- sqrt(mse)
cat("RMSE of GAM2: ", rmse)
```

```
## RMSE of GAM2:  0.7348539
```

## Comparing RMSE using 10-Fold Cross Validation for Linear Regression, GAM1, GAM2, KNN, and Lasso

```r
# Get data
wineData <- read.csv("~/Downloads/wine-quality-white-and-red.csv")

wineData$type <- as.factor(wineData$type)

# Split to training and testing data
set.seed(1)
train = sample(1:nrow(wineData), 0.7 * nrow(wineData))
train_data = wineData[train, ]
test_data = wineData[-train, ]

# Normalize function for KNN
normalize <- function(x) {
  return((x - mean(x)) / sd(x))
}


wine <- wineData

# Normalize numeric columns except quality
all_columns <- names(wine)
columns_to_normalize <- all_columns[all_columns != "quality" & sapply(wine, is.numeric)]
wine[columns_to_normalize] <- lapply(wine[columns_to_normalize], normalize)

# Convert factor
wine$type <- as.numeric(factor(wine$type))

# Setup cross-validation
k <- 10
folds <- sample(1:k, nrow(wine), replace = TRUE)

# Empty vectors to store RMSE
rmse_linear <- rmse_gam1 <- rmse_gam2 <- rmse_knn <- rmse_lasso <- numeric(k)

for (i in 1:k) {
  ktrain_data <- wine[folds != i, ]
  ktest_data <- wine[folds == i, ]

  # Linear Regression
  lm_model <- lm(quality ~ ., data = ktrain_data)
  linear_preds <- predict(lm_model, newdata = ktest_data)
  rmse_linear[i] <- sqrt(mean((linear_preds - ktest_data$quality)^2))

  # GAM Model 1 (smooth all variables)
  gam_model1 <- gam(quality ~
                    s(fixed.acidity) +
                    s(volatile.acidity) +
                    s(citric.acid) +
                    s(chlorides) +
                    s(residual.sugar) +
                    s(free.sulfur.dioxide) +
                    s(total.sulfur.dioxide) +
                    s(density) +
```

```r
                    s(sulphates) +
                    s(alcohol) + s(pH), data = ktrain_data, select=TRUE)
    gam_preds1 <- predict(gam_model1, newdata = ktest_data)
    rmse_gam1[i] <- sqrt(mean((gam_preds1 - ktest_data$quality)^2))

    # GAM Model 2 (smooth most variables)
    gam_model2 <- gam(quality ~
                        fixed.acidity +
                    volatile.acidity +
                    s(citric.acid) +
                    residual.sugar +
                    s(chlorides) +
                    free.sulfur.dioxide +
                    total.sulfur.dioxide +
                    density +
                    s(sulphates) +
                    s(alcohol) + s(pH),
                     data = ktrain_data)
    gam_preds2 <- predict(gam_model2, newdata = ktest_data)
    rmse_gam2[i] <- sqrt(mean((gam_preds2 - ktest_data$quality)^2))

  # KNN
  train_knn <- ktrain_data
  test_knn <- ktest_data
  train_knn$quality <- as.factor(train_knn$quality)
  test_knn$quality <- as.factor(test_knn$quality)

  training_X <- dplyr::select(train_knn, -quality)
  testing_X <- dplyr::select(test_knn, -quality)

  knn_preds <- knn.reg(train = training_X, test = testing_X, y = as.numeric(train_knn$quality), k = 10)
  rmse_knn[i] <- sqrt(mean((knn_preds - as.numeric(test_knn$quality))^2))

  # Lasso
  x_train <- as.matrix(dplyr::select(ktrain_data, -quality))
  y_train <- ktrain_data$quality
  x_test <- as.matrix(dplyr::select(ktest_data, -quality))

  lasso_cv <- cv.glmnet(x_train, y_train, alpha = 1)
  best_lambda <- lasso_cv$lambda.min
  lasso_preds <- predict(lasso_cv, s = best_lambda, newx = x_test)
  rmse_lasso[i] <- sqrt(mean((lasso_preds - ktest_data$quality)^2))
}

# Average RMSE across folds
results <- data.frame(
  Model = c("Linear", "GAM1", "GAM2", "KNN", "Lasso"),
  RMSE = c(mean(rmse_linear), mean(rmse_gam1), mean(rmse_gam2), mean(rmse_knn), mean(rmse_lasso))
)


print(results)
```

```
##     Model      RMSE
## 1 Linear 0.7336196
## 2   GAM1 0.7149347
## 3   GAM2 0.7291629
## 4    KNN 0.6903177
## 5  Lasso 0.7336021
```