

BACS HW (Week11)

106070038

2021-05-09

Question 1 Model fit is often determined by R^2 so let's dig into what this perspective of model fit is all about.

a. Let's dig into what regression is doing to compute model fit:

i. Plot Scenario 2, storing the returned points

```
plot_regr_rsqr <- function(points) {
  max_x <- 50
  if (nrow(points) == 0) {
    plot(NA, xlim=c(-5,max_x), ylim=c(-5,max_x), xlab="x", ylab="y")
    return()
  }
  plot(points, xlim=c(-5,max_x), ylim=c(-5,max_x), pch=19, cex=2, col="gray")
  if (nrow(points) < 2) return()

  mean_x <- mean(points$x)
  mean_y <- mean(points$y)
  segments(0, mean_y, max_x, mean_y, lwd=1, col="lightgray", lty="dotted")
  segments(mean_x, 0, mean_x, mean_y, lwd=1, col="lightgray", lty="dotted")
  regr <- lm(points$y ~ points$x)
  abline(regr, lwd=2, col="cornflowerblue")

  regr_summary <- summary(regr)
  ssr <- sum((regr$fitted.values - mean(points$y))^2)
  sse <- sum((points$y - regr$fitted.values)^2)
  sst <- sum((points$y - mean(points$y))^2)

  par(family="mono")
  legend("topleft", legend = c(
    paste(" Raw intercept: ", round(regr$coefficients[1], 2), "\n",
      "Raw slope      : ", round(regr$coefficients[2], 2), "\n",
      "Correlation    : ", round(cor(points$x, points$y), 2), "\n",
      "SSR           : ", round(ssr, 2), "\n",
      "SSE           : ", round(sse, 2), "\n",
      "SST           : ", round(sst, 2), "\n",
      "R-squared      : ", round(regr_summary$r.squared, 2))),
    bty="n")
  par(family="sans")
}
```

```

interactive_regression_rsq <- function(points=data.frame()) {
  cat("Click on the plot to create data points; hit [esc] to stop")
  repeat {
    plot_regr_rsq(points)
    click_loc <- locator(1)
    if (is.null(click_loc)) break
    if(nrow(points) == 0 ) {
      points <- data.frame(x=click_loc$x, y=click_loc$y)
    } else {
      points <- rbind(points, c(click_loc$x, click_loc$y))
    }
  }
  return(points)
}

# pts <- interactive_regression_rsq()
# save(pts, file='/Users/nkust/Desktop/2021Spring_Courses/BACS/HW9/pts.Rda')

```

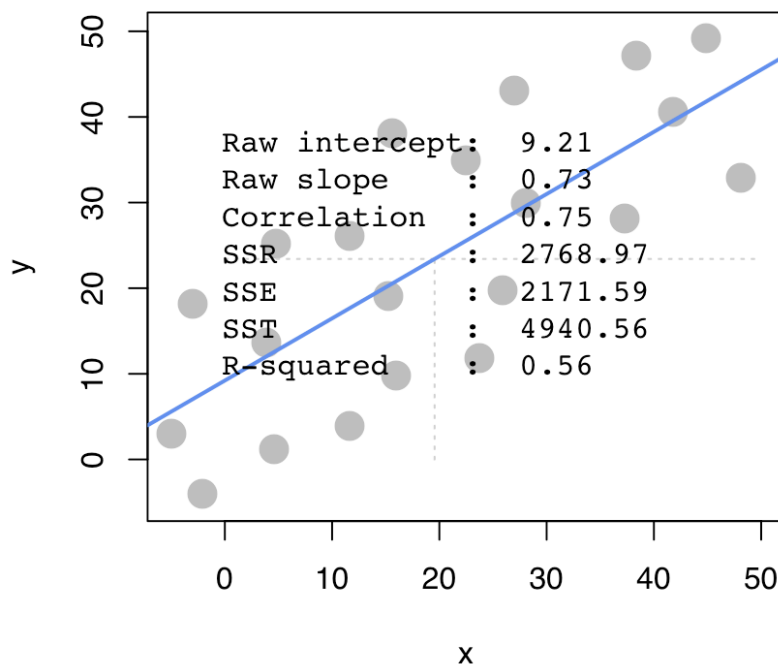


Figure 1: pts plot

- ii. Run a linear model of x and y points to confirm the R² value reported by the simulation

```
load('/Users/nkust/Desktop/2021Spring_Courses/BACS/HW9/pts.Rda')
```

```

regr <- lm(y ~ x, data=pts)
summary(regr)

##
## Call:
## lm(formula = y ~ x, data = pts)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -17.026  -8.449  -1.022   8.548  28.486
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  14.1505     3.9322   3.599  0.00179 **
## x              0.5431     0.1588   3.419  0.00272 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.64 on 20 degrees of freedom
## Multiple R-squared:  0.3689, Adjusted R-squared:  0.3374
## F-statistic: 11.69 on 1 and 20 DF,  p-value: 0.002717

```

- iii. Add line segments to the plot to show the regression residuals (errors)

```

y_hat <- regr$fitted.values
# segments(pts$x, pts$y, pts$x, y_hat, col="red", lty="dotted")

```

- iv. Use only `pts$x`, `pts$y`, `y_hat` and `mean(pts$y)` to compute SSE, SSR and SST, and verify R^2

```

SSE <- sum((y_hat - mean(pts$y))^2)
SSE

## [1] 1869.386

SSR <- sum((pts$y - y_hat)^2)
SSR

## [1] 3197.729

SST <- SSE + SSR
SST

## [1] 5067.114

R_square <- 1 - (SSR/(SSE + SSR))
R_square

```

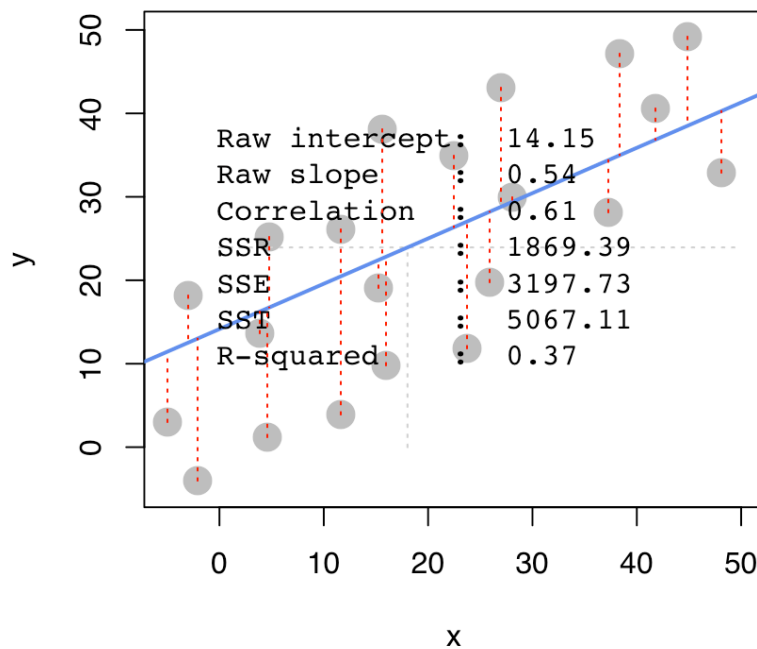


Figure 2: pts plot with line segments

```
## [1] 0.3689251
```

```
## verify
```

```
summary(regr)$r.square
```

```
## [1] 0.3689251
```

b. Comparing scenarios 1 and 2, which do we expect to have a stronger R^2 ?

- Ans: Scenario 1, because Scenario 1 is more intensive and close to the line than Scenario 2.

c. Comparing scenarios 3 and 4, which do we expect to have a stronger R^2 ?

- Ans: Scenario 3, because Scenario 3 is more intensive and close to the line than Scenario 4.

d. Comparing scenarios 1 and 2, which do we expect has bigger/smaller SSE, SSR, and SST? (do not compute SSE/SSR/SST here – just provide your intuition)

- Ans:
 - SSE: Scenario 1 < Scenario 2

- SSR: Scenario 1 > Scenario 2
- SST: Scenario 1 \approx Scenario 2

e. Comparing scenarios 3 and 4, which do we expect has bigger/smaller SSE, SSR, and SST? (do not compute SSE/SSR/SST here – just provide your intuition)

- Ans:
 - SSE: Scenario 3 < Scenario 4
 - SSR: Scenario 3 < Scenario 4
 - SST: Scenario 3 < Scenario 4

Question 2 We're going to take a look back at the early heady days of global car manufacturing, when American, Japanese, and European cars competed to rule the world. Take a look at the data set in file auto-data.txt. We are interested in explaining what kind of cars have higher fuel efficiency (mpg).

```
auto <- read.table("auto-data.txt", header=FALSE, na.strings = "?")
names(auto) <- c("mpg", "cylinders", "displacement", "horsepower", "weight", "acceleration", "model_year")
```

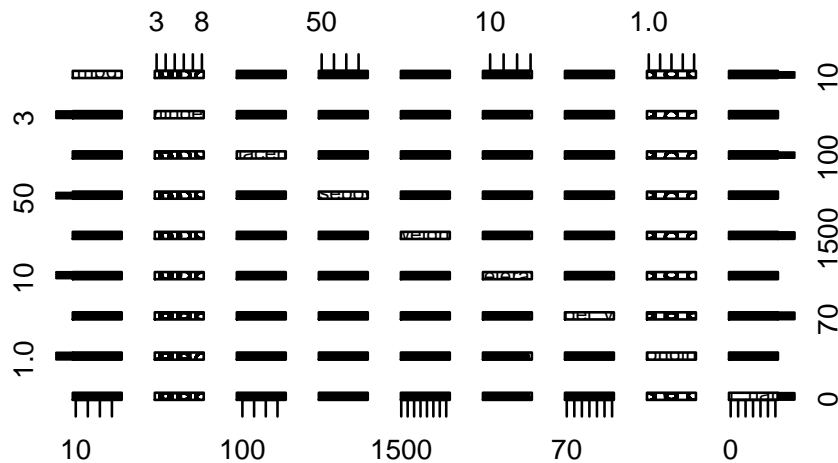
a.

i. Visualize the data in any way you feel relevant

```
summary(auto)
```

```
##      mpg      cylinders      displacement      horsepower      weight
##  Min.   : 9.00   Min.   :3.000   Min.   : 68.0   Min.   : 46.0   Min.   :1613
## 1st Qu.:17.50   1st Qu.:4.000   1st Qu.:104.2   1st Qu.: 75.0   1st Qu.:2224
## Median :23.00   Median :4.000   Median :148.5   Median : 93.5   Median :2804
## Mean   :23.51   Mean   :5.455   Mean   :193.4   Mean   :104.5   Mean   :2970
## 3rd Qu.:29.00   3rd Qu.:8.000   3rd Qu.:262.0   3rd Qu.:126.0   3rd Qu.:3608
## Max.   :46.60   Max.   :8.000   Max.   :455.0   Max.   :230.0   Max.   :5140
##
##      acceleration      model_year      origin      car_name
##  Min.   : 8.00   Min.   :70.00   Min.   :1.000   Length:398
## 1st Qu.:13.82   1st Qu.:73.00   1st Qu.:1.000   Class :character
## Median :15.50   Median :76.00   Median :1.000   Mode  :character
## Mean   :15.57   Mean   :76.01   Mean   :1.573
## 3rd Qu.:17.18   3rd Qu.:79.00   3rd Qu.:2.000
## Max.   :24.80   Max.   :82.00   Max.   :3.000
##
```

```
plot(auto)
```



- ii. Report a correlation table of all variables, rounding to two decimal places

```
cor <- cor(auto[-9], use="pairwise.complete.obs")
round(cor,2)
```

```
##           mpg cylinders displacement horsepower weight acceleration
## mpg      1.00   -0.78      -0.80      -0.78  -0.83      0.42
## cylinders -0.78    1.00       0.95       0.84   0.90     -0.51
## displacement -0.80  0.95       1.00       0.90   0.93     -0.54
## horsepower  -0.78  0.84       0.90       1.00   0.86     -0.69
## weight     -0.83  0.90       0.93       0.86   1.00     -0.42
## acceleration 0.42 -0.51     -0.54     -0.69 -0.42      1.00
## model_year  0.58 -0.35     -0.37     -0.42 -0.31      0.29
## origin     0.56 -0.56     -0.61     -0.46 -0.58      0.21
##           model_year origin
## mpg           0.58   0.56
## cylinders     -0.35 -0.56
## displacement  -0.37 -0.61
## horsepower    -0.42 -0.46
## weight        -0.31 -0.58
## acceleration  0.29  0.21
## model_year    1.00  0.18
## origin        0.18  1.00
```

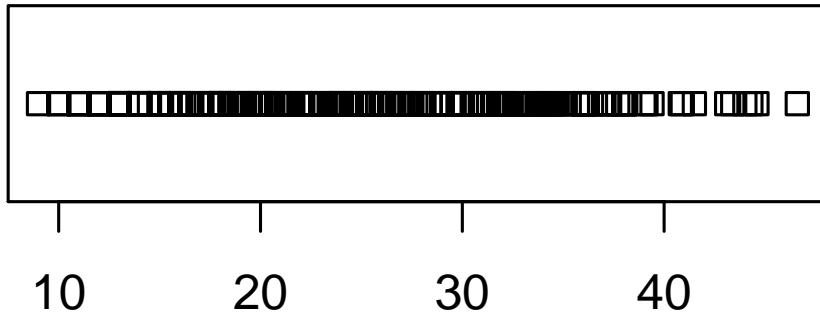
- iii. From the visualizations and correlations, which variables seem to relate to mpg?

```
cor(auto[1], auto[-9], use="pairwise.complete.obs")
```

```
##           mpg cylinders displacement horsepower weight acceleration model_year
## mpg      1 -0.7753963  -0.8042028 -0.7784268 -0.8317409   0.4202889  0.5792671
##           origin
## mpg      0.5634504
```

- Ans: displacement(-0.8042028) and weight(-0.8317409)
- iv. Which relationships might not be linear? (don't worry about linearity for rest of this HW)

```
plot(auto[1])
```



- Ans: Based on the correlation plot, most of the relationships with model_year are not linear.

- v. Are there any pairs of independent variables that are highly correlated ($r > 0.7$)?

```
summary(auto[-9])
```

```
##      mpg      cylinders  displacement  horsepower      weight
##  Min.   : 9.00   Min.   :3.000   Min.   : 68.0   Min.   : 46.0   Min.   :1613
## 1st Qu.:17.50   1st Qu.:4.000   1st Qu.:104.2   1st Qu.: 75.0   1st Qu.:2224
## Median :23.00   Median :4.000   Median :148.5   Median : 93.5   Median :2804
## Mean   :23.51   Mean   :5.455   Mean   :193.4   Mean   :104.5   Mean   :2970
## 3rd Qu.:29.00   3rd Qu.:8.000   3rd Qu.:262.0   3rd Qu.:126.0   3rd Qu.:3608
## Max.   :46.60   Max.   :8.000   Max.   :455.0   Max.   :230.0   Max.   :5140
##
##      acceleration  model_year      origin
##  Min.   : 8.00   Min.   :70.00   Min.   :1.000
## 1st Qu.:13.82   1st Qu.:73.00   1st Qu.:1.000
## Median :15.50   Median :76.00   Median :1.000
## Mean   :15.57   Mean   :76.01   Mean   :1.573
## 3rd Qu.:17.18   3rd Qu.:79.00   3rd Qu.:2.000
## Max.   :24.80   Max.   :82.00   Max.   :3.000
##
```

- Ans:

b. Let's create a linear regression model where mpg is dependent upon all other suitable variables (Note: origin is categorical with three levels, so use `factor(origin)` in `lm(...)` to split it into two dummy variables)

- i. Which independent variables have a 'significant' relationship with mpg at 1% significance?
- ii. Looking at the coefficients, is it possible to determine which independent variables are the most effective at increasing mpg? If so, which ones, and if not, why not? (hint: units!)

c.