

Lab 3

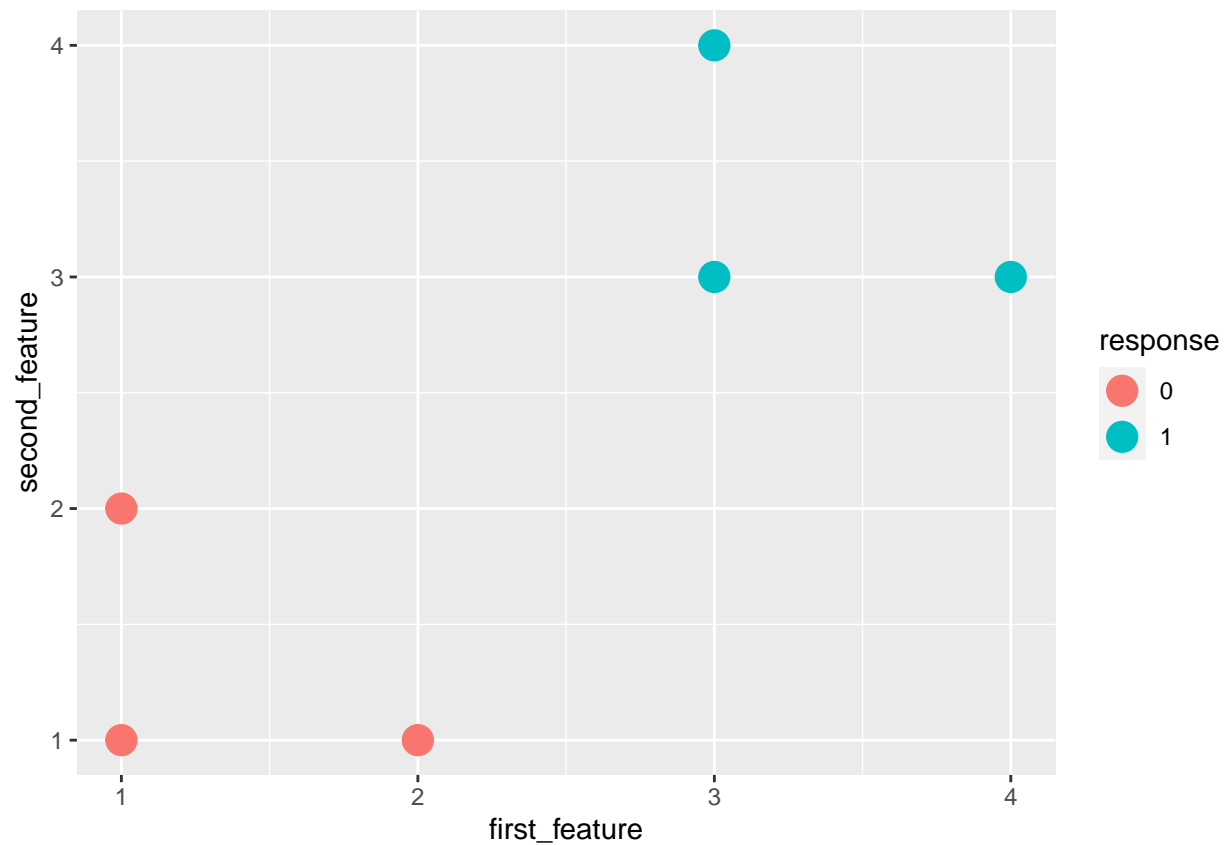
Alexis Montes

11:59PM March 4, 2021

Support Vector Machine vs. Perceptron

We recreate the data from the previous lab and visualize it:

```
pacman::p_load(ggplot2)
Xy_simple = data.frame(
  response = factor(c(0, 0, 0, 1, 1, 1)), #nominal
  first_feature = c(1, 1, 2, 3, 3, 4), #continuous
  second_feature = c(1, 2, 1, 3, 4, 3) #continuous
)
simple_viz_obj = ggplot(Xy_simple, aes(x = first_feature, y = second_feature, color = response)) +
  geom_point(size = 5)
simple_viz_obj
```

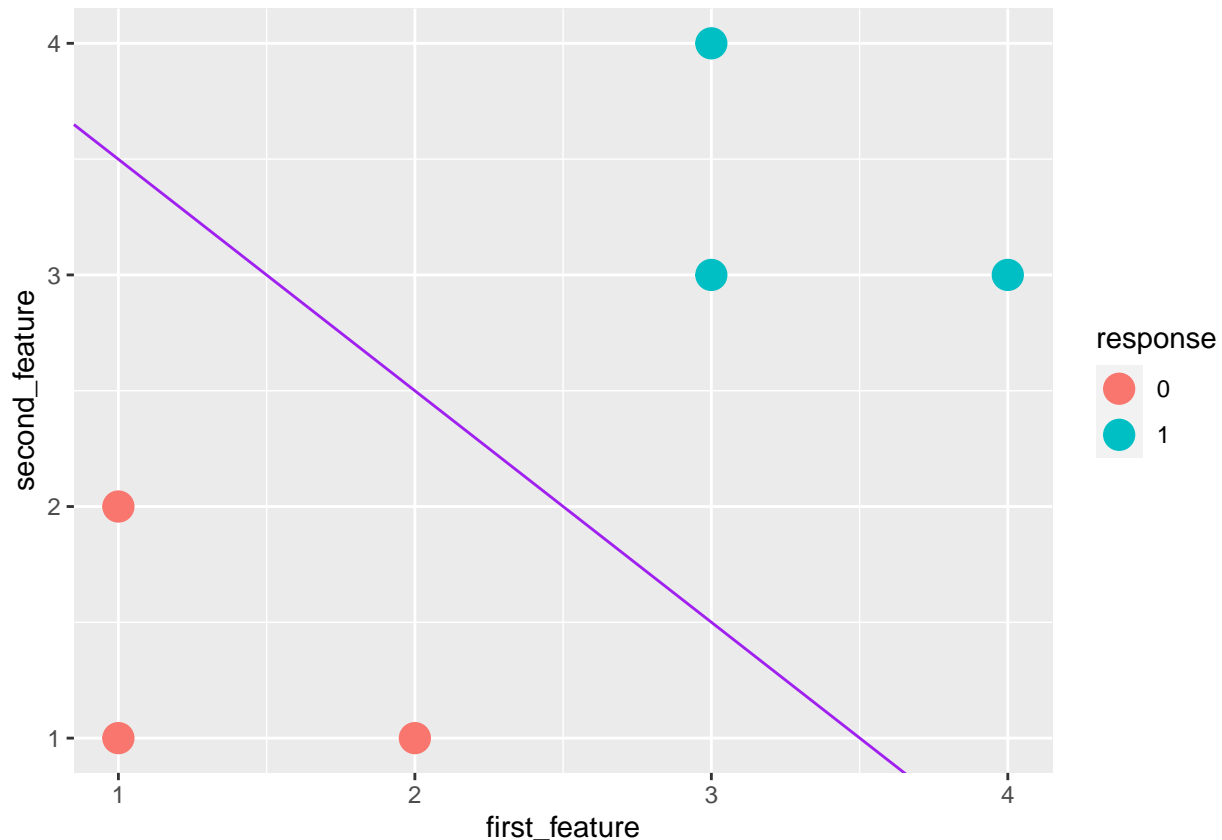


Use the `e1071` package to fit an SVM model to the simple data. Use a formula to create the model, pass in the data frame, set kernel to be `linear` for the linear SVM and don't scale the covariates. Call the model object `svm_model`. Otherwise the remaining code won't work.

```
pacman::p_load(e1071)
svm_model = svm(
  as.matrix(Xy_simple[, 2 :3]),
  Xy_simple$response,
  kernel = "linear",
  scale = FALSE
)
```

and then use the following code to visualize the line in purple:

```
w_vec_simple_svm = c(
  svm_model$rho, #the b term
  -t(svm_model$coefs) %*% cbind(Xy_simple$first_feature, Xy_simple$second_feature)[svm_model$index, ] #
)
simple_svm_line = geom_abline(
  intercept = -w_vec_simple_svm[1] / w_vec_simple_svm[3],
  slope = -w_vec_simple_svm[2] / w_vec_simple_svm[3],
  color = "purple"
)
simple_viz_obj + simple_svm_line
```



Source the `perceptron_learning_algorithm` function from lab 2. Then run the following to fit the perceptron and plot its line in orange with the SVM's line:

```

#' Perceptron learning algorithm
#'
#' The perceptron fits the best line for the passed data
#'
#' @param Xinput      Training data as a matrix
#' @param y_binary    Training data as a binary matrix
#' @param MAX_ITER    maximum number of iterations for the function
#' @param w           weights data as vector
#'
#' @return            The computed final parameter (weight) as a vector of length p + 1
perceptron_learning_algorithm = function(Xinput, y_binary, MAX_ITER = 1000, w = NULL){
  num_rows <- nrow(Xinput)

  Xinput <- cbind(rep(1, num_rows), Xinput)

  num_cols <- ncol(Xinput)

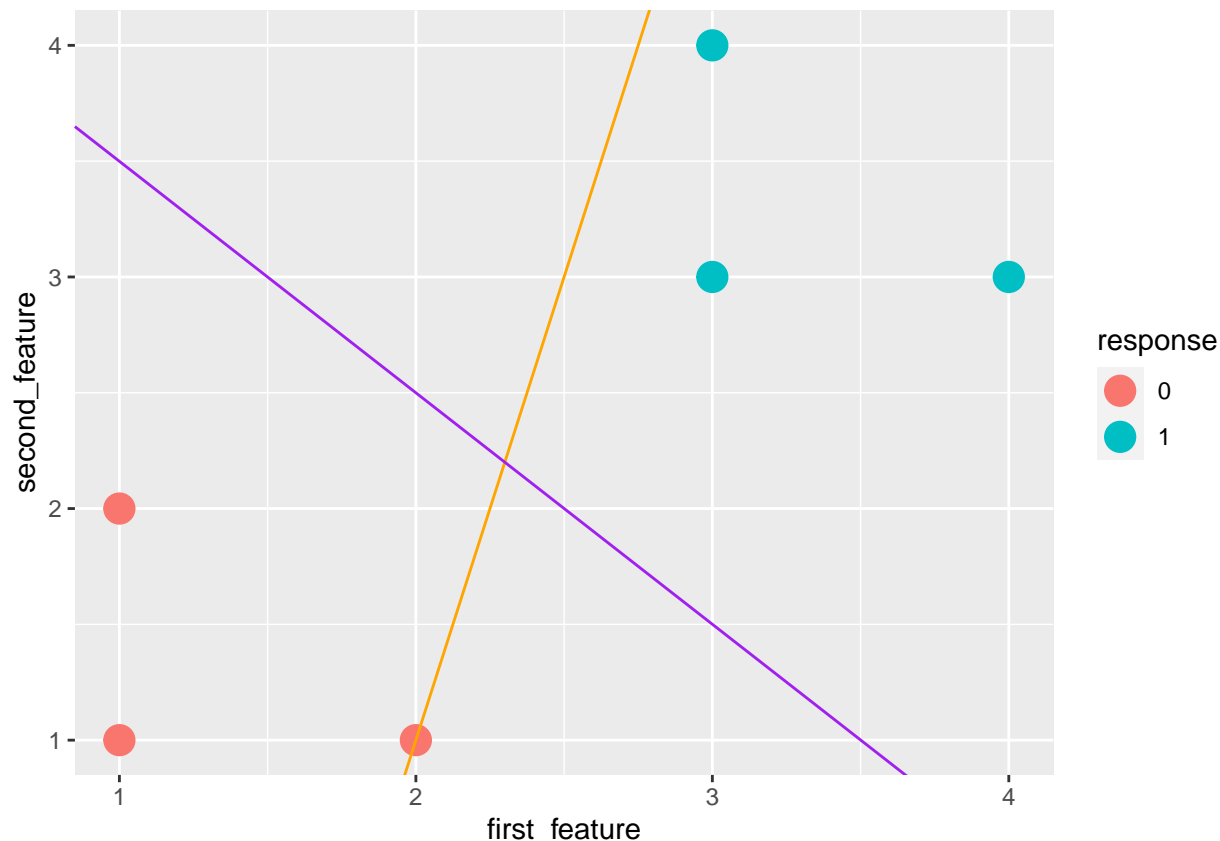
  if (is.null(w)) {
    w <- rep(0, num_cols)
  }

  for (i in 1:MAX_ITER) {
    for (j in 1:num_rows) {
      X <- Xinput[j, ]
      e <- y_binary[j] - ifelse(sum(w %*% X) > 0, 1, 0)
      for (k in 1:num_cols) {
        w[k] <- w[k] + e * X[k]
      }
    }
  }
  return(w)
}

w_vec_simple_per = perceptron_learning_algorithm(
  cbind(Xy_simple$first_feature, Xy_simple$second_feature),
  as.numeric(Xy_simple$response == 1)
)

simple_perceptron_line = geom_abline(
  intercept = -w_vec_simple_per[1] / w_vec_simple_per[3],
  slope = -w_vec_simple_per[2] / w_vec_simple_per[3],
  color = "orange")
simple_viz_obj + simple_perceptron_line + simple_svm_line

```



Is this SVM line a better fit than the perceptron?

TO-DO

Now write pseudocode for your own implementation of the linear support vector machine algorithm using the Vapnik objective function we discussed.

Note there are differences between this spec and the perceptron learning algorithm spec in question #1. You should figure out a way to respect the `MAX_ITER` argument value.

```
## Support Vector Machine
#
## This function implements the hinge-loss + maximum margin linear support vector machine algorithm of
#
## @param Xinput      The training data features as an n x p matrix.
## @param y_binary    The training data responses as a vector of length n consisting of only 0's and 1's.
## @param MAX_ITER    The maximum number of iterations the algorithm performs. Defaults to 5000.
## @param lambda      A scalar hyperparameter trading off margin of the hyperplane versus average hinge
#
## The default value is 1.
## @return            The computed final parameter (weight) as a vector of length p + 1
linear_svm_learning_algorithm = function(Xinput, y_binary, MAX_ITER = 5000, lambda = 0.1){

  # Initialize a w vector
  # for ( iter in 1 : MAX_ITER){
  #   for (i in 1 : nrow(Xinput)){
  #     Create a function that will help pick the maximum margin for the data
  #   }
  # }
}
```

```
#}
}
```

If you are enrolled in 342W the following is extra credit but if you're enrolled in 650, the following is required. Write the actual code. You may want to take a look at the `optimx` package. You can feel free to define another function (a “private” function) in this chunk if you wish. R has a way to create public and private functions, but I believe you need to create a package to do that (beyond the scope of this course).

```
#' This function implements the hinge-loss + maximum margin linear support vector machine algorithm of
#'
#' @param Xinput      The training data features as an n x p matrix.
#' @param y_binary    The training data responses as a vector of length n consisting of only 0's and 1's.
#' @param MAX_ITER    The maximum number of iterations the algorithm performs. Defaults to 5000.
#' @param lambda      A scalar hyperparameter trading off margin of the hyperplane versus average hinge
#'
#' The default value is 1.
#' @return            The computed final parameter (weight) as a vector of length p + 1
#linear_sum_learning_algorithm = function(Xinput, y_binary, MAX_ITER = 5000, lambda = 0.1){
  #TO-DO
#}
```

If you wrote code (the extra credit), run your function using the defaults and plot it in brown vis-a-vis the previous model's line:

Is this the same as what the `e1071` implementation returned? Why or why not?

TO-DO

We now move on to simple linear modeling using the ordinary least squares algorithm.

Let's quickly recreate the sample data set from practice lecture 7:

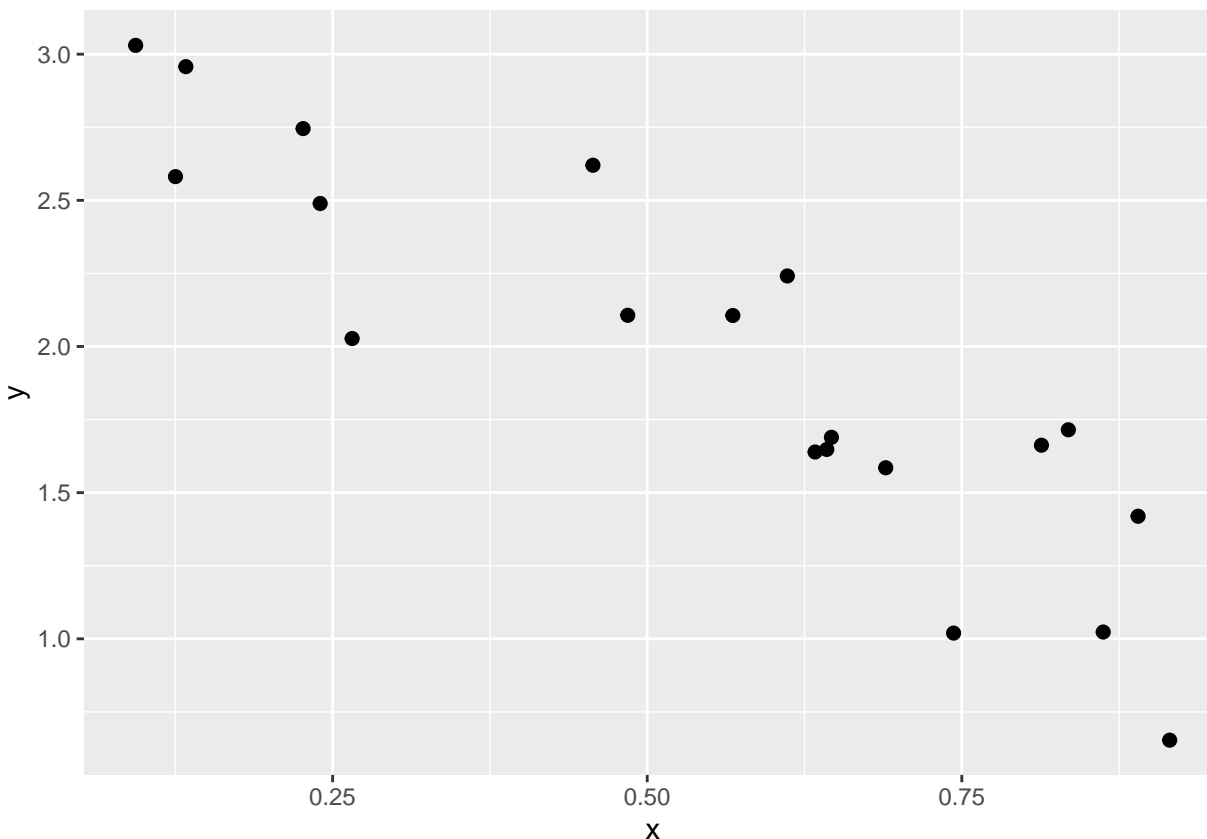
```
n = 20
x = runif(n)
beta_0 = 3
beta_1 = -2
```

Compute $h^*(x)$ as `h_star_x`, then draw $\epsilon \sim N(0, 0.33^2)$ as `epsilon`, then compute `y`.

```
h_star_x = beta_0 + beta_1 * x
epsilon = rnorm(n, 0, 0.33)
y = h_star_x + epsilon
```

Graph the data by running the following chunk:

```
pacman::p_load(ggplot2)
simple_df = data.frame(x = x, y = y)
simple_viz_obj = ggplot(simple_df, aes(x, y)) +
  geom_point(size = 2)
simple_viz_obj
```



Does this make sense given the values of β_0 and β_1 ?

Write a function `my_simple_ols` that takes in a vector `x` and vector `y` and returns a list that contains the `b_0` (intercept), `b_1` (slope), `yhat` (the predictions), `e` (the residuals), `SSE`, `SST`, `MSE`, `RMSE` and `Rsqr` (for the R-squared metric). Internally, you can only use the functions `sum` and `length` and other basic arithmetic operations. You should throw errors if the inputs are non-numeric or not the same length. You should also name the class of the return value `my_simple_ols_obj` by using the `class` function as a setter. No need to create ROxygen documentation here.

```
my_simple_ols = function(x, y){
  n = length(y)
  x_bar = sum(x)/length(x)
  y_bar = sum(y)/n
  if( length(x) != n){
    stop("x and y need to be the same length")
  }
  if(class(x) != 'numeric' && class(x) != 'integer' ){
    stop("x needs to be numeric")
  }
  if(class(y) != 'numeric' && class(y) != 'integer' ){
    stop("y needs to be numeric")
  }
  if(n < 2){
    stop("n must be more than 2")
  }

  b_1 = (sum(x*y)-n*x_bar*y_bar )/ (sum(x^2) - n*x_bar^2)
```

```

b_0 = y_bar - b_1 * x_bar
yhat = b_0 + b_1 * x
e = y - yhat
SSE = sum(e^2)
SST = sum((y-y_bar)^2)
MSE = SSE/ (n-2)
RMSE = sqrt(MSE)
Rsqr = 1 - SSE/SST

model = list(b_0 = b_0, b_1 = b_1, yhat = yhat, e = e, SSE = SSE, SST = SST, MSE = MSE, RMSE = RMSE, Rsqr = Rsqr)
class(model) = "my_simple_ols_obj"
model
}

```

Verify your computations are correct for the vectors `x` and `y` from the first chunk using the `lm` function in R:

```

lm_mod = lm(y~x)
my_simple_ols_mod = my_simple_ols(x,y)
#run the tests to ensure the function is up to spec
pacman::p_load(testthat)
expect_equal(my_simple_ols_mod$b_0, as.numeric(coef(lm_mod)[1]), tol = 1e-4)
expect_equal(my_simple_ols_mod$b_1, as.numeric(coef(lm_mod)[2]), tol = 1e-4)
expect_equal(my_simple_ols_mod$RMSE, summary(lm_mod)$sigma, tol = 1e-4)
expect_equal(my_simple_ols_mod$Rsqr, summary(lm_mod)$r.squared, tol = 1e-4)

```

Verify that the average of the residuals is 0 using the `expect_equal`. Hint: use the syntax above.

```
expect_equal(mean(my_simple_ols_mod$e), 0)
```

Create the `X` matrix for this data example. Make sure it has the correct dimension.

```
X = cbind(1,x)
```

Use the `model.matrix` function to compute the matrix `X` and verify it is the same as your manual construction.

```
model.matrix(~x)
```

```
##      (Intercept)          x
## 1             1 0.45685871
## 2             1 0.86237604
## 3             1 0.91518752
## 4             1 0.74346947
## 5             1 0.24010313
## 6             1 0.09346238
## 7             1 0.64269430
## 8             1 0.63342423
## 9             1 0.13337182
## 10            1 0.48453827
## 11            1 0.22654017

```

```
## 12          1 0.64639572
## 13          1 0.56803784
## 14          1 0.83465631
## 15          1 0.26550069
## 16          1 0.61136520
## 17          1 0.81328516
## 18          1 0.12512080
## 19          1 0.68954017
## 20          1 0.89006101
## attr("assign")
## [1] 0 1
```

Create a prediction method `g` that takes in a vector `x_star` and `my_simple_ols_obj`, an object of type `my_simple_ols_obj` and predicts `y` values for each entry in `x_star`.

```
g = function(my_simple_ols_obj, x_star){
  my_simple_ols_obj$b_0 + my_simple_ols_obj$b_1 * x_star
}
```

Use this function to verify that when predicting for the average `x`, you get the average `y`.

```
expect_equal(g(my_simple_ols_mod, mean(x)), mean(y))
```

In class we spoke about error due to ignorance, misspecification error and estimation error. Show that as `n` grows, estimation error shrinks. Let us define an error metric that is the difference between b_0 and b_1 and β_0 and β_1 . How about $h = \|b - \beta\|^2$ where the quantities are now the vectors of size two. Show as `n` increases, this shrinks.

```
beta_0 = 3
beta_1 = -2
beta = c(beta_0, beta_1)
ns = 10^(1:6)
errors_in_b = array(NA, length(ns))
for (i in 1 : length(ns)) {
  n = ns[i]
  x = runif(n)
  h_star_x = beta_0 + beta_1 * x
  epsilon = rnorm(n, mean = 0, sd = 0.33)
  y = h_star_x + epsilon

  mod = my_simple_ols(x,y)
  b = c(mod$b_0, mod$b_1)
  errors_in_b[i] = sum((beta - b)^2)
}
log(errors_in_b,10)
```

```
## [1] -0.8822387 -1.7746995 -3.1833021 -3.8994769 -4.4715088 -7.3270888
```

```
errors_in_b
```

```
## [1] 1.311479e-01 1.679966e-02 6.556890e-04 1.260443e-04 3.376690e-05
## [6] 4.708810e-08
```


We are now going to repeat one of the first linear model building exercises in history — that of Sir Francis Galton in 1886. First load up package `HistData`.

```
pacman::p_load(HistData)
```

In it, there is a dataset called `Galton`. Load it up.

```
data(Galton)
```

You now should have a data frame in your workspace called `Galton`. Summarize this data frame and write a few sentences about what you see. Make sure you report n , p and a bit about what the columns represent and how the data was measured. See the help file `?Galton`. p is 1 and n is 928 the number of observations

```
pacman::p_load(skimr)
skim(Galton)
```

Table 1: Data summary

Name	Galton
Number of rows	928
Number of columns	2
Column type frequency:	
numeric	2
Group variables	None

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
parent	0	1	68.31	1.79	64.0	67.5	68.5	69.5	73.0	
child	0	1	68.09	2.52	61.7	66.2	68.2	70.2	73.7	

TO-DO

Find the average height (include both parents and children in this computation).

```
avg_height = mean(c(Galton$parent, Galton$child))
```

If you were to use the null model, what would the RMSE be of this model be?

```
n = nrow(Galton)
SST = sum((Galton$child - mean(Galton$child))^2)
sqrt(SST/(n-1))
```

```
## [1] 2.517941
```

Note that in Math 241 you learned that the sample average is an estimate of the “mean”, the population expected value of height. We will call the average the “mean” going forward since it is probably correct to the nearest tenth of an inch with this amount of data.

Run a linear model attempting to explain the childrens' height using the parents' height. Use `lm` and use the R formula notation. Compute and report b_0 , b_1 , RMSE and R^2 .

```
mod = lm(child~parent,Galton)
b_0 = coef(mod)[1]
b_1 = coef(mod)[2]
summary(mod)$sigma
```

```
## [1] 2.238547
```

```
summary(mod)$r.squared
```

```
## [1] 0.2104629
```

Interpret all four quantities: b_0 , b_1 , RMSE and R^2 . Use the correct units of these metrics in your answer.

TO-DO

How good is this model? How well does it predict? Discuss.

TO-DO

It is reasonable to assume that parents and their children have the same height? Explain why this is reasonable using basic biology and common sense.

TO-DO

If they were to have the same height and any differences were just random noise with expectation 0, what would the values of β_0 and β_1 be?

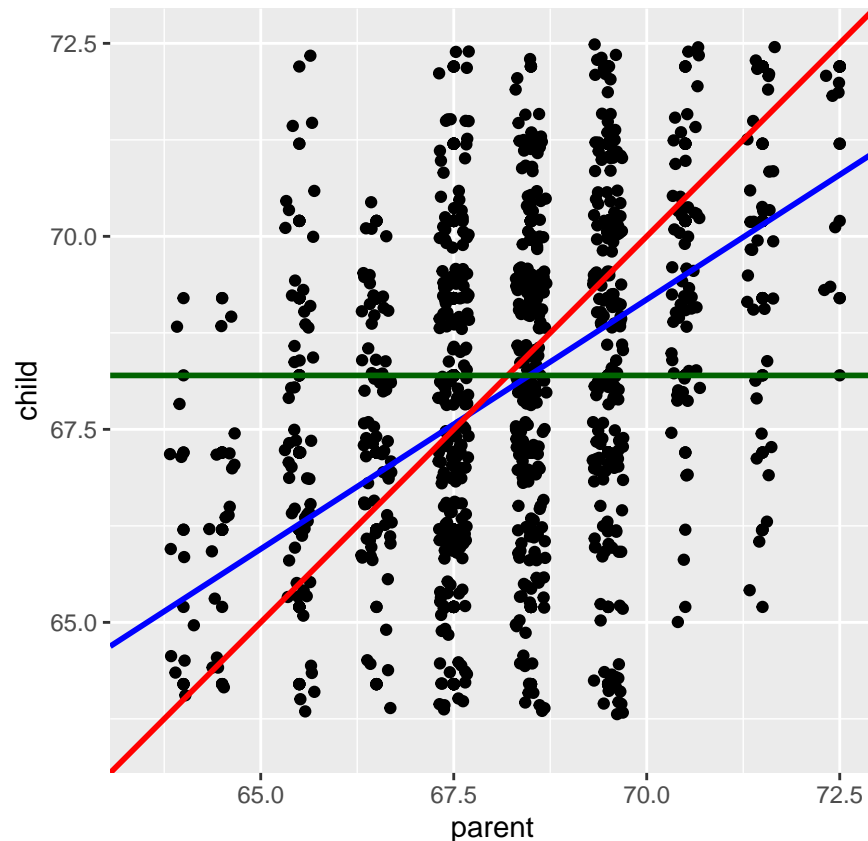
TO-DO

Let's plot (a) the data in \mathbb{D} as black dots, (b) your least squares line defined by b_0 and b_1 in blue, (c) the theoretical line β_0 and β_1 if the parent-child height equality held in red and (d) the mean height in green.

```
pacman::p_load(ggplot2)
ggplot(Galton, aes(x = parent, y = child)) +
  geom_point() +
  geom_jitter() +
  geom_abline(intercept = b_0, slope = b_1, color = "blue", size = 1) +
  geom_abline(intercept = 0, slope = 1, color = "red", size = 1) +
  geom_abline(intercept = avg_height, slope = 0, color = "darkgreen", size = 1) +
  xlim(63.5, 72.5) +
  ylim(63.5, 72.5) +
  coord_equal(ratio = 1)
```

```
## Warning: Removed 76 rows containing missing values (geom_point).
```

```
## Warning: Removed 87 rows containing missing values (geom_point).
```



Fill in the following sentence:

TO-DO: Children of short parents became tall on average and children of tall parents became shorter on average.

Why did Galton call it “Regression towards mediocrity in hereditary stature” which was later shortened to “regression to the mean”?

TO-DO: The average height of the children regressed to the mean or average of the height of the parents. The children of taller parents ended up being shorter, averaging towards the mean height. And the children of shorter parents ended up being taller, averaging towards the mean about.

Why should this effect be real?

This effect possibly occurred due to sampling error. Since results are taken from children that are tall or short, the two will average out towards the mean. When taking a similar sample size from two opposing values they will always regress because their values would average each other out.

You now have unlocked the mystery. Why is it that when modeling with y continuous, everyone calls it “regression”? Write a better, more descriptive and appropriate name for building predictive models with y continuous.

TO-DO It is called a regression because of Galton’s findings that the average heights of short children and tall children move towards the average height, showing the “regression” of the heights.

You can now clear the workspace. Create a dataset \mathbb{D} which we call Xy such that the linear model as R^2 about 50% and RMSE approximately 1.

```
x = c(10, 7, 7, 8, 8, 5)
y = c(2, 4, 3, 5, 3, 5)
Xy = data.frame(x = x, y = y)
```

Create a dataset \mathbb{D} which we call `Xy` such that the linear model as R^2 about 0% but `x`, `y` are clearly associated.

```
x = c(4, 2, 4, 5, 7, 8)
y = c(2, 4, 2, 4, 2, 4)
Xy = data.frame(x = x, y = y)
```

Extra credit: create a dataset \mathbb{D} and a model that can give you R^2 arbitrarily close to 1 i.e. approximately 1 - epsilon but RMSE arbitrarily high i.e. approximately M .

```
epsilon = 0.01
M = 1000
#TO-DO
```

Write a function `my_ols` that takes in `X`, a matrix with with `p` columns representing the feature measurements for each of the `n` units, a vector of `n` responses `y` and returns a list that contains the `b`, the $p + 1$ -sized column vector of OLS coefficients, `yhat` (the vector of `n` predictions), `e` (the vector of `n` residuals), `df` for degrees of freedom of the model, `SSE`, `SST`, `MSE`, `RMSE` and `Rsquared` (for the R-squared metric). Internally, you cannot use `lm` or any other package; it must be done manually. You should throw errors if the inputs are non-numeric or not the same length. Or if `X` is not otherwise suitable. You should also name the class of the return value `my_ols` by using the `class` function as a setter. No need to create ROxygen documentation here.

```
my_ols = function(X, y){
  n = nrow(X)
  p = ncol(X) + 1

  if (nrow(X) != nrow(y)) {
    stop("X and y must have the same number of rows")
  }
  if (class(X) != "numeric" & class(X) != "integer") {
    stop("x needs to be numeric")
  }
  if (class(y) != "numeric" & class(y) != "integer") {
    stop("y needs to be numeric")
  }
  if (n <= 2) {
    stop("Dimensions of n must be greater than 2")
  }

  X = as.matrix(cbind(1, X[, 1:p]))
  colnames(X)[1] = "(intercept)"
  XtX = t(X) %*% X

  tryCatch({
    XtXinv = solve(XtX)
  }, error = function(e) {
    stop("Matrix is not invertible or does not have full rank")
  })

  b = XtXinv %*% t(X) %*% y
  yhat = X %*% b
  e = y - yhat
  df = ...
```

```

SSE = t(e) %*% e
s_sq_y = var(y )
SST = (n -1) * s_sq_y
MSE = 1 / (nrow(X) - ncol(X)) * SSE
RMSE = sqrt(MSE)
Rsqr = 1 - SSE / SST

model = list("b"=b, "yhat"=yhat, "e"=e, "SSE"=SSE, "SST"=SST, "MSE"=MSE, "RMSE"=RMSE, "Rsqr"=Rsqr)

class(model) = "my_ols"

model
}

```

Verify that the OLS coefficients for the **Type** of cars in the cars dataset gives you the same results as we did in class (i.e. the \bar{y} 's within group).

#TO-DO

Create a prediction method **g** that takes in a vector **x_star** and the dataset \mathbb{D} i.e. **X** and **y** and returns the OLS predictions. Let **X** be a matrix with with **p** columns representing the feature measurements for each of the **n** units

```

g = function(x_star, X, y){
  model = my_ols(X, y)

  x_star %*% model$b
}

```