Lab 3 **Brendan Gubbins** 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 3 second\_feature response 0 first\_feature 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) Xy\_simple\_feature\_matrix = as.matrix(Xy\_simple[, 2 : 3])  $svm_model = svm($ x = Xy\_simple\_feature\_matrix, # not sure why I needed this to remove error formula = Xy\_simple\_feature\_matrix, data = Xy\_simple\$response, kernel = "linear", scale = FALSE  $\#svm\_model = svm(Xy\_simple\_feature\_matrix, Xy\_simple\$response, kernel = "linear", cost = (2 * n * lambda)^-1, sca$ 1e = 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, ] # the other 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 3 second\_feature response 2 first\_feature 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: # sourcing my perceptron function from lab02 perceptron\_learning\_algorithm = function(Xinput, y\_binary, MAX\_ITER = 1000, w = NULL){ Xinput = as.matrix(cbind(1, Xinput[,, drop = FALSE])) n = nrow(Xinput)p = ncol(Xinput) w = rep(0, p)for (i in 1 : MAX\_ITER) { **for** (j **in** 1 : n) { x\_j = Xinput[j,]  $y_hat = ifelse(sum(x_j * w) > 0, 1, 0)$ for (k in 1 : p) {  $w[k] = w[k] + (y\_binary[j] - y\_hat) * x\_j[k]$ 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 3 second\_feature response first\_feature Is this SVM line a better fit than the perceptron? The purple SVM line is clearly far superior to the perceptron, because it cleanly splits the data. The orange perceptron line is touching the data at (2,1).Now write pseuocode 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 Vladimir V apnik (1963). #' @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 loss. The default value is 1. #' @return The computed final parameter (weight) as a vector of length p + 1linear\_svm\_learning\_algorithm = function(Xinput, y\_binary, MAX\_ITER = 5000, lambda = 0.1){ # for 1 to MAX\_ITER # initialize w and SHE variables # for 1 to nrow(Xinput) # set w to satisfy such that ||w|| satisfies  $(y\_binary\_i - 1/2)(w * Xinput[i, ] - b) >= 1/2$  $# min_w = min(w, min_w)$ # compute sum of the hinge errors (SHE) by summing  $\max(0, 1/2 - (y_binary_i - 1/2)(w * Xinput[i,] - b))$ # optimize for argmin w,b  $(1/nrow(Xinput) * SHE + lambda * ||w||^2)$ # return w } 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 Vladimir V apnik (1963). # ' #' @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 loss. #' 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){ #T0-D0 } 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: svm\_model\_weights = linear\_svm\_learning\_algorithm(X\_simple\_feature\_matrix, y\_binary) my\_svm\_line = geom\_abline( intercept = svm\_model\_weights[1] / svm\_model\_weights[3],#NOTE: negative sign removed from intercept argument slope = -svm\_model\_weights[2] / svm\_model\_weights[3], color = "brown") simple\_viz\_obj + my\_svm\_line Is this the same as what the e1071 implementation returned? Why or why not? 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 = 20x = 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, mean = 0, sd = 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 3.0 -2.5 -2.0 -1.5 -1.0 -0.25 0.75 0.00 0.50 Does this make sense given the values of  $\beta_0$  and  $\beta_1$ ? beta 0 is 3 and beta 1 is -2. The graph makes sense because if you fit a line on these data points it will be about a -2 slope (beta 1) and the intercept is about 3 (beta\_0) 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 Rsq (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 my\_simple\_ols = function(x, y){ n = length(y)**if** (n != length(x)) { stop("x and y must be same length.") if (class(x) != 'numeric' & class(y) != 'integer') { stop("x needs to be numeric.") if (class(y) != 'numeric' & class(x) != 'integer') { stop("y needs to be numeric.") } **if** (n <= 2) { stop("n must be more than 2.") }  $x_bar = sum(x)/n$  $y_bar = sum(y)/n$  $b_1 = (sum(x^y) - n * x_bar * y_bar)/(sum(x^2) - n * x_bar^2)$  $b_0 = y_bar - b_1 * x_bar$  $yhat = b_0 + b_1 * x$ e = y - yhat $SSE = sum(e^2)$  $SST = sum((y - y_bar)^2)$ MSE = SSE / (n-2)RMSE = sqrt(MSE)Rsq = 1 - SSE / SST $model = list(b_0 = b_0, b_1 = b_1, yhat = yhat, e = e, SSE = SSE, SST = SST, MSE = MSE, RMSE = RMSE, Rsq = Rsq)$ class(model) = "my\_simple\_ols\_obj" model # return Verify your computations are correct for the vectors x and y from the first chunk using the 1m function in R:  $lm_mod = lm(y \sim x)$ 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\$Rsq, 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. residuals = mean(my\_simple\_ols\_mod\$e) # expect\_equal(residuals, 0) # not good in real world, residuals may go to 0 on R expect\_equal(residuals, 0, tol = 1e-4) Create the X matrix for this data example. Make sure it has the correct dimension. X = cbind(1, x) # cbind makes a matrixΧ ## [1,] 1 0.35116932 ## [2,] 1 0.95597895 ## [3,] 1 0.31373789 ## [4,] 1 0.36909007 ## [5,] 1 0.07867264 ## [6,] 1 0.02017286 ## [7,] 1 0.15689565 ## [8,] 1 0.10244965 ## [9,] 1 0.00994406 ## [10,] 1 0.25882570 ## [11,] 1 0.39219961 ## [12,] 1 0.33639667 ## [13,] 1 0.92495885 ## [14,] 1 0.02316466 ## [15,] 1 0.03334887 ## [16,] 1 0.60214046 ## [17,] 1 0.70046498 ## [18,] 1 0.95212643 ## [19,] 1 0.64266330 ## [20,] 1 0.05788750 Use the model.matrix function to compute the matrix X and verify it is the same as your manual construction.  $model.matrix( \sim x)$ ## (Intercept) ## 1 1 0.35116932 ## 2 1 0.95597895 ## 3 1 0.31373789 ## 4 1 0.36909007 ## 5 1 0.07867264 ## 6 1 0.02017286 ## 7 1 0.15689565 ## 8 1 0.10244965 ## 9 1 0.00994406 ## 10 1 0.25882570 ## 11 1 0.39219961 ## 12 1 0.33639667 ## 13 1 0.92495885 ## 14 1 0.02316466 ## 15 1 0.03334887 ## 16 1 0.60214046 ## 17 1 0.70046498 ## 18 1 0.95212643 ## 19 1 0.64266330 ## 20 1 0.05788750 ## 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){ y\_star = my\_simple\_ols\_obj\$b\_0 + my\_simple\_ols\_obj\$b\_1 \* x\_star y\_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  $b_2$  and  $b_3$ . How about  $h=||b-\beta||^2$  where the quantities are now the vectors of size two. Show as n increases, this shrinks.  $beta_0 = 3$  $beta_1 = -2$  $beta = c(beta_0, beta_1)$  $ns = 10^{(1:8)}$ error\_in\_b = array(NA, length(ns)) # errors in b not being beta 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)  $y = h_star_x + epsilon$  $mod = my\_simple\_ols(x,y)$  $b = c(mod\$b\_0, mod\$b\_1)$  $error_in_b[i] = sum((beta - b)^2)$ log(error\_in\_b, 10) ## [1] -2.227607 -3.095099 -3.659925 -3.579467 -5.115667 -6.068457 -9.742476 ## [8] -6.824372 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)

Galton

928

2

2

69.5

70.2

p0

64.0

61.7

p25

67.5

66.2

p50

68.5

68.2

None

p100 hist

73.0

73.7

Data summary

Number of rows

Number of columns

Column type frequency:

Name

numeric

Group variables

skim\_variable

as a balancing effect.

n = nrow(Galton)

## [1] 2.517941

coef(mod)

## (Intercept)

## 23.9415302

 $b_0 = coef(mod)[1]$  $b_1 = coef(mod)[2]$ 

## [1] 2.238547

## [1] 0.2104629

sense.

summary(mod)\$r.squared

inheriting DNA from their parents.

pacman::p\_load(ggplot2)

xlim(63.5, 72.5) +ylim(63.5, 72.5) +coord\_equal(ratio = 1)

72.5 -

70.0 -

Fill in the following sentence:

Why should this effect be real?

rm(list = ls())

 $model = lm(y \sim x)$ 

## [1] 0.501566

## [1] 1.236338

x = c(1, 4, 8, 12)y = c(4,1,2,3)

 $model = lm(y \sim x)$ 

## [1] 0.01818182

summary(model)\$sigma

x = c(7.5, 11, 9, 9, 1)y = c(2, 5, 3, 5, 2)

summary(model)\$r.squared

Xy = data.frame(x = x, y = y)

Xy = data.frame(x = x, y = y)

summary(model)\$r.squared

height.

geom\_point() + geom\_jitter() +

report  $b_0$ ,  $b_1$ , RMSE and  $R^2$ .

mod = lm(child ~ parent, Galton)

parent

child

Variable type: numeric

n missing

0

0

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

avg\_height = mean(c(Galton\$parent, Galton\$child))

 $SST = sum((Galton\$child - mean(Galton\$child))^2)$ sqrt(SST/(n-1)) # not n-2, only fitting one ybar

parent

summary(mod)\$sigma # +- 2.23, +- 4.5 95% of the time

21% which is fairly low, this means that the model does not fit the variance well.

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

parent height. Any parent height produces the same exact child height.

child height equality held in red and (d) the mean height in green.

ggplot(Galton, aes(x = parent, y = child)) +

0.6462906

complete\_rate

1

If you were predicting child height from parent height and you were using the null mode, what would the RMSE of the null model be?

average the "mean" going forward since it is probably correct to the nearest tenth of an inch with this amount of data.

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

R^2 is low (21% of the variance explained), so the model is not accounting for the data variance.

geom\_abline(intercept = b\_0, slope = b\_1, color = "blue", size = 1) +

geom\_abline(intercept = avg\_height, slope = 0, color = "darkgreen", size = 1) +

67.5

producing taller and taller children with no end. This effect keeps humans generally of similar size on average.

be called "historical sampling," which takes historical data as samples and spits out predictions based on that data.

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.

descriptive and appropriate name for building predictive models with y continuous.

parent

70.0

Children of short parents became taller than expected on average and children of tall parents became shorter than expected on average.

Because the children of very tall parents were shorter than expected, and therefore closer to the mean height, and the children of very short parents were taller than expected, closer to the mean height. These children of very tall parents/very short parents "regressed" towards the mean

This effect should be real because short parents would keep producing smaller and smaller children with no limit, and tall parents would keep

It's called regression because of Galton's coinage of "Regression" to the mean. A more appropriate description/name for predictive models could

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.

You now have unlocked the mystery. Why is it that when modeling with y continuous, everyone calls it "regression"? Write a better, more

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

72.5

geom\_abline(intercept = 0, slope = 1, color = "red", size = 1) +

## Warning: Removed 76 rows containing missing values (geom\_point).

## Warning: Removed 84 rows containing missing values (geom\_point).

because being off by 4.46 inches in height 95% of the time is significant (a range of between 5'4" to 6'+).

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

Run a linear model attempting to explain the childrens' height using the parents' height. Use 1m and use the R formula notation. Compute and

 $b_0$  is the intercept in inches of the child's height. In this case, if the average of father and mother heights are 0 (absurd), the child will be 23.9~ inches tall.  $b_1$  is the increase in average child height per 1 increase in average height of father and mother (child gets 0.65~ inches taller per 1 inch of average parent height). The RMSE is 2.23 $\sim$ , meaning that 95% of the time, there is a +- of 4.46 $\sim$  inches on the child's predicted height.  $R^2$  is

b 0 intercept/if parent height is 0, child height is 23.94 inch. b 1 increase in child (.64) height per average height of parent. RMSE tells the range.

It's a pretty bad model because from the RMSE, you can tell somebody's height 95% of the time from a +- of about 4.46 inches. This is not good

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

It's reasonable for the most part to assume that children will be the same height as their parents i.e taller parents will produce taller children and shorter parents will produce shorter children on average. It would be unreasonable to assume that children height is randomized, because they are

If they were to have the same height and any differences were just random noise with expectation 0, what would the values of  $\beta_0$  and  $\beta_1$  be?

beta 0 would be the intercept at 0 and beta 1 would be 1. This means there would be a 1:1 change when determining child height based on

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  $\beta_0$  and  $\beta_1$  if the parent-

68.31

68.09

The number of rows n = 928, so there are 928 parents heights averaged, as well as a corresponding child height. The two columns represent parents and child, p = 2. The mean averaged height of parents is 68.3 inches, and the mean height of children is 68.1 inches. The lowest height found in the 0th percentile is 61.7 inches (child), and the largest height found in the 100th percentile is 73.7 inches (child). Based on the histogram, the vast majority of height is found in the 25th-75th percentile. As described in the dataset, all female children had their heights multiplied by 1.08

1.79

2.52

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 = 0.01M = 1000x = c(5, 90, 30, 15, 2, 50, 44, 1020, 999)y = c(50, 1000, 700, 899, 2000, 11, 2, 10000, 7802)Xy = data.frame(x = x, y = y) $model = lm(y \sim x)$ summary(model)\$r.squared ## [1] 0.9449758 summary(model)\$sigma ## [1] 934.52 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 Rsq (for the R-squared metric). Internally, you cannot use 1m 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){ X = cbind(1, X)n = length(y)p = ncol(X)**if** (n != nrow(X)) { stop("X # rows and y must be same length.") } if (class(X[1,]) != 'numeric' & class(y) != 'integer') { stop("X needs to be numeric.") if (class(y) != 'numeric' & class(X[1,]) != 'integer') { stop("y needs to be numeric.") } **if** (n <= p + 1) { stop("n must be more than p+1.") } df = p + 1Xt = t(X)XtX = Xt %\*% XXtXinv = solve(XtX)b = XtXinv %\*% Xt %\*% y $y_bar = sum(y)/n$ yhat = X %\*% be = y - yhatSSE = t(e) % % e $SST = sum((y - y_bar)^2)$ MSE = SSE /(n - (p + 1))RMSE = sqrt(MSE)Rsq = 1 - SSE / SSTmodel = list(b = b, yhat = yhat, e = e, df = df, SSE = SSE, SST = SST, MSE = MSE, RMSE = RMSE, Rsq = Rsq) class(model) = "my\_ols" model # return 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 ybar's within group). cars = MASS::Cars93 X\_mm = model.matrix(~Type, cars)[,2:6]  $mod = lm(Price \sim Type, cars)$ X = my\_ols(X\_mm, cars\$Price) X\$b ## [,1] ## 18.212500 ## TypeLarge 6.087500 ## TypeMidsize 9.005682 ## TypeSmall -8.045833 ## TypeSporty 1.180357 ## TypeVan 0.887500 mod ## ## Call: ## lm(formula = Price ~ Type, data = cars) ## Coefficients: ## (Intercept) TypeLarge TypeMidsize TypeSmall TypeSporty 18.2125 6.0875 9.0057 -8.0458 1.1804

TypeVan

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){

 $y_hat = c(1, x_star) %*% b$ 

g(X\_mm[5,], X\_mm, cars\$Price)

[,1]

## [1,] 27.21818

predict(mod, cars[5,])

5

## 27.21818

 $b = my_ols(X, y)$ \$b

y\_hat

0.8875