

# Stat 432 Homework 7

Assigned: Oct 7, 2024; Due: 11:59 PM CT, Oct 17, 2024

- Instruction
- Question 1: SVM on Hand Written Digit Data (55 points)
- Question 2: SVM with Kernel Trick (45 points)

## Instruction

**Please remove this section when submitting your homework.**

Students are encouraged to work together on homework and/or utilize advanced AI tools. However, **sharing, copying, or providing any part of a homework solution or code to others** is an infraction of the University's rules on Academic Integrity (<https://studentcode.illinois.edu/article1/part4/1-401/>). Any violation will be punished as severely as possible. Final submissions must be uploaded to Gradescope (<https://www.gradescope.com/courses/570816>). No email or hard copy will be accepted. For **late submission policy and grading rubrics** (<https://teazrq.github.io/stat432/syllabus.html>), please refer to the course website.

- You are required to submit the rendered file `HWx_yourNetID.pdf`. For example, `HW01_rqzhu.pdf`. Please note that this must be a `.pdf` file. `.html` format **cannot** be accepted. Make all of your R code chunks visible for grading.
- Include your Name and NetID in the report.
- If you use this file or the example homework `.Rmd` file as a template, be sure to **remove this instruction** section.
- Make sure that you **set seed** properly so that the results can be replicated if needed.
- For some questions, there will be restrictions on what packages/functions you can use. Please read the requirements carefully. As long as the question does not specify such restrictions, you can use anything.
- **When using AI tools**, you are encouraged to document your comment on your experience with AI tools especially when it's difficult for them to grasp the idea of the question.
- **On random seed and reproducibility**: Make sure the version of your R is  $\geq 4.0.0$ . This will ensure your random seed generation is the same as everyone else. Please note that updating the R version may require you to reinstall all of your packages.

## Question 1: SVM on Hand Written Digit Data (55 points)

We will again use the MNIST dataset. We will use the first 2400 observations of it:

```

# inputs to download file
fileLocation <- "https://pjreddie.com/media/files/mnist_train.csv"
numRowsToDownload <- 2400
localFileName <- paste0("mnist_first", numRowsToDownload, ".RData")

# download the data and add column names
mnist2400 <- read.csv(fileLocation, nrow = numRowsToDownload)
numColsMnist <- dim(mnist2400)[2]
colnames(mnist2400) <- c("Digit", paste("Pixel", seq(1:(numColsMnist - 1))), sep = "")

# save file
# in the future we can read in from the local copy instead of having to redownload
save(mnist2400, file = localFileName)

# you can load the data with the following code
#load(file = localFileName)

```

- a. [15 pts] Since a standard SVM can only be used for binary classification problems, let's fit SVM on digits 4 and 5. Complete the following tasks.
  - Use digits 4 and 5 in the first 1200 observations as training data and those in the remaining part with digits 4 and 5 as testing data.
  - Fit a linear SVM on the training data using the `e1071` package. Set the cost parameter  $C = 1$ .
  - You will possibly encounter two issues: first, this might be slow (unless your computer is very powerful); second, the package will complain about some pixels being problematic (zero variance). Hence, reducing the number of variables by removing pixels with low variances is probably a good idea. Perform a marginal screening of variance on the pixels and select the top 250 Pixels with the highest marginal variance.
  - Redo your SVM model with the pixels you have selected. Report the training and testing classification errors.
- b. [15 pts] Some researchers might be interested in knowing what pixels are more important in distinguishing the two digits. One way to do this is to extract the coefficients of the (linear) SVM model (they are fairly comparable in our case since all the variables have the same range). Keep in mind that the coefficients are those  $\beta$  parameter used to define the direction of the separation line, and they can be recovered from the solution of the Lagrangian. Complete the following tasks.
  - Extract the coefficients of the linear SVM model you have fitted in part 1. State the mathematical formula of how these coefficients are recovered using the solution of the Lagrangian.
  - Find the top 30 pixels with the largest absolute coefficients.
  - Refit the SVM using just these 30 pixels. Report the training and testing classification errors.
- c. [15 pts] Perform a logistic regression with elastic net penalty ( $\alpha = 0.5$ ) on the training data. Start with the 250 pixels you have used in part a). You do not need to select the best  $\lambda$  value using cross-validation. Instead, select the model with just 30 variables in the solution path (what is this? you can refer to our lecture note on Lasso). What is the  $\lambda$  value corresponding to this model? Extract the pixels being selected by your elastic net model. Do these pixels overlap with the ones selected by the SVM model in part b)? Comment on your findings.
- d. [10 pts] Compare the two 30-variable models you obtained from part b) and c). Use the area under the ROC curve (AUC) on the testing data as the performance metric.

## Question 2: SVM with Kernel Trick (45 points)

This problem involves the `OJ` data set which is part of the `ISLR2` package. We create a training set containing a random sample of 800 observations, and a test set containing the remaining observations. In the dataset, `Purchase` variable is the output variable and it indicates whether a customer purchased Citrus Hill or Minute Maid Orange Juice. For the details of the dataset you can refer to its help file.

```
library(ISLR2)
data("OJ")
set.seed(7)
id=sample(nrow(OJ),800)
train=OJ[id,]
test=OJ[-id,]
```

- [15 pts]\*\* Fit a (linear) support vector machine by using `svm` function to the training data using `cost = 0.01` and using all the input variables. Provide the training and test errors.
- [15 pts]\*\* Use the `tune()` function to select an optimal cost, `C` in the set of  $\{0.01, 0.1, 1, 2, 5, 7, 10\}$ . Compute the training and test errors using the best value for cost.
- [15 pts]\*\* Repeat parts 1 and 2 using a support vector machine with `radial` and `polynomial` (with degree 2) kernel. Use the default value for `gamma` in the `radial` kernel. Comment on your results from parts b and c.