

# MATH 342W / 642 / RM 742 Spring 2024 HW #5

Professor Adam Kapelner

Due 11:59PM May 15

(this document last updated 10:18am on Friday 3<sup>rd</sup> May, 2024)

## Instructions and Philosophy

The path to success in this class is to do many problems. Unlike other courses, exclusively doing reading(s) will not help. Coming to lecture is akin to watching workout videos; thinking about and solving problems on your own is the actual “working out.” Feel free to “work out” with others; **I want you to work on this in groups.**

Reading is still *required*. You should be googling and reading about all the concepts introduced in class online. This is your responsibility to supplement in-class with your own readings.

The problems below are color coded: **green** problems are considered *easy* and marked “[easy]”; **yellow** problems are considered *intermediate* and marked “[harder]”, **red** problems are considered *difficult* and marked “[difficult]” and **purple** problems are extra credit. The *easy* problems are intended to be “giveaways” if you went to class. Do as much as you can of the others; I expect you to at least attempt the *difficult* problems.

This homework is worth 100 points but the point distribution will not be determined until after the due date. See syllabus for the policy on late homework.

Up to 7 points are given as a bonus if the homework is typed using L<sup>A</sup>T<sub>E</sub>X. Links to installing L<sup>A</sup>T<sub>E</sub>X and program for compiling L<sup>A</sup>T<sub>E</sub>X is found on the syllabus. You are encouraged to use **overleaf.com**. If you are handing in homework this way, read the comments in the code; there are two lines to comment out and you should replace my name with yours and write your section. The easiest way to use overleaf is to copy the raw text from hwxx.tex and preamble.tex into two new overleaf tex files with the same name. If you are asked to make drawings, you can take a picture of your handwritten drawing and insert them as figures or leave space using the “\vspace” command and draw them in after printing or attach them stapled.

The document is available with spaces for you to write your answers. If not using L<sup>A</sup>T<sub>E</sub>X, print this document and write in your answers. I do not accept homeworks which are *not* on this printout. Keep this first page printed for your records.

NAME: \_\_\_\_\_

## Problem 1

These are some questions related to probability estimation modeling and asymmetric cost modeling.

- (a) [easy] Why is logistic regression an example of a “generalized linear model” (glm)?
  
  
  
  
  
  
  
  
  
  
- (b) [easy] What is  $\mathcal{H}_{pr}$  for the probability estimation algorithm that employs the linear model in the covariates with logistic link function?
  
  
  
  
  
  
  
  
  
  
- (c) [easy] If logistic regression predicts 3.1415 for a new  $\mathbf{x}_*$ , what is the probability estimate that  $y = 1$  for this  $\mathbf{x}_*$ ?
  
  
  
  
  
  
  
  
  
  
- (d) [harder] What is  $\mathcal{H}_{pr}$  for the probability estimation algorithm that employs the linear model in the covariates with cloglog link function?
  
  
  
  
  
  
  
  
  
  
- (e) [difficult] Generalize linear probability estimation to the case where  $\mathcal{Y} = \{C_1, C_2, C_3\}$ . Use the logistic link function like in logistic regression. Write down the objective function that you would numerically maximize. This objective function is one that is

argmax'd over the parameters (you define what these parameters are — that is part of the question).

Once you get the answer you can see how this easily goes to  $K > 3$  response categories. The algorithm for general  $K$  is known as “multinomial logistic regression”, “polytomous LR”, “multiclass LR”, “softmax regression”, “multinomial logit” (mlogit), the “maximum entropy” (MaxEnt) classifier, and the “conditional maximum entropy model”. You can inflate your resume with lots of jazz by doing this one question!

- (f) [easy] Graph a canonical ROC and label the axes. In your drawing estimate AUC. Explain very clearly what is measured by the  $x$  axis and the  $y$  axis.

- (g) [easy] Pick one point on your ROC curve from the previous question. Explain a situation why you would employ this model.
- (h) [harder] Graph a canonical DET curve and label the axes. Explain very clearly what is measured by the  $x$  axis and the  $y$  axis. Make sure the DET curve's intersections with the axes is correct.
- (i) [easy] Pick one point on your DET curve from the previous question. Explain a situation why you would employ this model.
- (j) [difficult] [MA] The line of random guessing on the ROC curve is the diagonal line with slope one extending from the origin. What is the corresponding line of random guessing in the DET curve? This is not easy...

## Problem 2

These are some questions related to bias-variance decomposition. Assume the two assumptions from the notes about the random variable model that produces the  $\delta$  values, the error due to ignorance.

- (a) [easy] Write down (do not derive) the decomposition of MSE for a given  $\mathbf{x}_*$  where  $\mathbb{D}$  is assumed fixed but the response associated with  $\mathbf{x}_*$  is assumed random.
  
- (b) [easy] Write down (do not derive) the decomposition of MSE for a given  $\mathbf{x}_*$  where the responses in  $\mathbb{D}$  is random but the  $\mathbf{X}$  matrix is assumed fixed and the response associated with  $\mathbf{x}_*$  is assumed random like previously.
  
- (c) [easy] Write down (do not derive) the decomposition of MSE for general predictions of a phenomenon where all quantities are considered random.
  
- (d) [difficult] Why is it in (a) there is only a “bias” but no “variance” term? Why did the additional source of randomness in (b) spawn the variance term, a new source of error?

- (e) [harder] A high bias / low variance algorithm is underfit or overfit?
- (f) [harder] A low bias / high variance algorithm is underfit or overfit?
- (g) [harder] Explain why bagging reduces MSE for “free” regardless of the algorithm employed.
  
- (h) [harder] Explain why RF reduces MSE atop bagging  $M$  trees and specifically mention the target that it attacks in the MSE decomposition formula and why it’s able to reduce that target.
  
  
  
  
  
  
  
  
  
  
- (i) [difficult] When can RF lose to bagging  $M$  trees? Hint: think hyperparameter choice.

### Problem 3

These are some questions related to missingness.

- (a) [easy] [MA] What are the three missing data mechanisms? Provide an example when each occurs (i.e., a real world situation). We didn't really cover this in class so I'm making it a MA question only. This concept will NOT be on the exam.
  
  
  
  
  
  
  
  
  
  
- (b) [easy] Why is listwise-deletion a *terrible* idea to employ in your  $\mathbb{D}$  when doing supervised learning?
  
  
  
  
  
  
  
  
  
  
- (c) [easy] Why is it good practice to augment  $\mathbb{D}$  to include missingness dummies? In other words, why would this increase oos predictive accuracy?
  
  
  
  
  
  
  
  
  
  
- (d) [easy] To impute missing values in  $\mathbb{D}$ , what is a good default strategy and why?

## Problem 4

These are some questions related to gradient boosting. The final gradient boosted model after  $M$  iterations is denoted  $G_M$  which can be written in a number of equivalent ways (see below). The  $g_t$ 's denote constituent models and the  $G_t$ 's denote partial sums of the constituent models up to iteration number  $t$ . The constituent models are “steps in functional steps” which have a step size  $\eta$  and a direction component denoted  $\tilde{g}_t$ . The directional component is the base learner  $\mathcal{A}$  fit to the negative gradient of the objective function  $L$  which measures how close the current predictions are to the real values of the responses:

$$\begin{aligned} G_M &= G_{M-1} + g_M \\ &= g_0 + g_1 + \dots + g_M \\ &= g_0 + \eta \tilde{g}_1 + \dots + \eta \tilde{g}_M \\ &= g_0 + \eta \mathcal{A}(\langle \mathbf{X}, -\nabla L(\mathbf{y}, \hat{\mathbf{y}}_1) \rangle, \mathcal{H}) + \dots + \eta \mathcal{A}(\langle \mathbf{X}, -\nabla L(\mathbf{y}, \hat{\mathbf{y}}_M) \rangle, \mathcal{H}) \\ &= g_0 + \eta \mathcal{A}(\langle \mathbf{X}, -\nabla L(\mathbf{y}, g_1(\mathbf{X})) \rangle, \mathcal{H}) + \dots + \eta \mathcal{A}(\langle \mathbf{X}, -\nabla L(\mathbf{y}, g_M(\mathbf{X})) \rangle, \mathcal{H}) \end{aligned}$$

- (a) [easy] From a perspective of only multivariable calculus, explain gradient descent and why it's a good idea to find the minimum inputs for an objective function  $L$  (in English).

- (b) [easy] Write the mathematical steps of gradient boosting for supervised learning below. Use  $L$  for the objective function to keep the procedure general. Use notation found in the problem header.



(c) [easy] For regression, what is  $g_0(\mathbf{x})$ ?

(d) [easy] For probability estimation for binary response, what is  $g_0(\mathbf{x})$ ?

(e) [harder] What are all the hyperparameters of gradient boosting? There are more than just two.

(f) [easy] For regression, rederive the negative gradient of the objective function  $L$ .

(g) [easy] For probability estimation for binary response, rederive the negative gradient of the objective function  $L$ .

- (h) [difficult] For probability estimation for binary response scenarios, what is the unit of the output  $G_M(\mathbf{x}_*)$ ?
- (i) [easy] For the base learner algorithm  $\mathcal{A}$ , why is it a good idea to use shallow CART (which is the recommended default)?
- (j) [difficult] For the base learner algorithm  $\mathcal{A}$ , why is it a bad idea to use deep CART?
- (k) [difficult] For the base learner algorithm  $\mathcal{A}$ , why is it a bad idea to use OLS for regression (or logistic regression for probability estimation for binary response)?
- (l) [difficult] If  $M$  is very, very large, what is the risk in using gradient boosting even using shallow CART as the base learner (the recommended default)?
- (m) [difficult] If  $\eta$  is very, very large but  $M$  reasonably correctly chosen, what is the risk in using gradient boosting even using shallow CART as the base learner (the recommended default)?