# MATH 342W / 650.4 Spring 2021 Homework #3

#### Professor Adam Kapelner

Due 11:59PM Wednesday, April 7, 2021 by email

(this document last updated  $10:36\,\mathrm{pm}$  on Wednesday  $17^\mathrm{th}$  March, 2021)

#### 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 LATEX. Links to instaling LATEX and program for compiling LATEX 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 LATEX, 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 questions about Silver's book, chapters 3-6. For all parts in this question, answer using notation from class (i.e.  $t, f, g, h^*, \delta, \epsilon, e, t, z_1, \ldots, z_t, \mathbb{D}, \mathcal{H}, \mathcal{A}, \mathcal{X}, \mathcal{Y}, X, y, n, p, x_{-1}, \ldots, x_{-p}, x_{1}, \ldots, x_{n}$ , etc.

(a) [difficult] Chapter 4 is all about predicting weather. Broadly speaking, what is the problem with weather predictions? Make sure you use the framework and notation from class. This is not an easy question and we will discuss in class. Do your best.

(b) [easy] Why does the weatherman lie about the chance of rain? And where should you go if you want honest forecasts?

(c) [difficult] Chapter 5 is all about predicting earthquakes. Broadly speaking, what is the problem with earthquake predictions? It is *not* the same as the problem of predicting weather. Read page 162 a few times. Make sure you use the framework and notation from class.

(d)	[easy] Silver has quite a whimsical explanation of overfitting on page 163 but it is really educational! What is the nonsense predictor in the model he describes?
(e)	[easy] John von Neumann was credited with saying that "with four parameters I can fit an elephant and with five I can make him wiggle his trunk". What did he mean by that and what is the message to you, the budding data scientist?
(f)	[difficult] Chapter 6 is all about predicting unemployment, an index of macroeconomic performance of a country. Broadly speaking, what is the problem with unemployment predictions? It is not the same as the problem of predicting weather or earthquakes Make sure you use the framework and notation from class.
(g)	[E.C.] Many times in this chapter Silver says something on the order of "you need to have theories about how things function in order to make good predictions." Do you agree? Discuss.

## Problem 2

These are questions related to the concept of orthogonal projection, QR decomposition and its relationship with least squares linear modeling.

(a) [difficult] [MA] Prove that if a square matrix is both symmetric and idempotent then it must be an orthogonal projection matrix.

(b) [easy] Prove that  $I_n$  is an orthogonal projection matrix  $\forall n.$ 

(c) [easy] What subspace does  $I_n$  project onto?

(d) [easy] Consider least squares linear regression using a design matrix X with rank p+1. What are the degrees of freedom in the resulting model? What does this mean?

(e) [easy] If you are orthogonally projecting the vector  $\boldsymbol{y}$  onto the column space of X which is of rank p+1, derive the formula for  $\operatorname{Proj}_{\operatorname{colsp}[X]}[\boldsymbol{y}]$ . Is this the same as in OLS?

(f) [difficult] We saw that the perceptron is an *iterative algorithm*. This means that it goes through multiple iterations in order to converge to a closer and closer  $\boldsymbol{w}$ . Why not do the same with linear least squares regression? Consider the following. Regress  $\boldsymbol{y}$  using  $\boldsymbol{X}$  to get  $\hat{\boldsymbol{y}}$ . This generates residuals  $\boldsymbol{e}$  (the leftover piece of  $\boldsymbol{y}$  that wasn't explained by the regression's fit,  $\hat{\boldsymbol{y}}$ ). Now try again! Regress  $\boldsymbol{e}$  using  $\boldsymbol{X}$  and then get new residuals  $\boldsymbol{e}_{new}$ . Would  $\boldsymbol{e}_{new}$  be closer to  $\boldsymbol{0}_n$  than the first  $\boldsymbol{e}$ ? That is, wouldn't this yield a better model on iteration #2? Yes/no and explain.

(g) [harder] Prove that  $Q^{\top} = Q^{-1}$  where Q is an orthonormal matrix such that colsp  $[Q] = \operatorname{colsp}[X]$  and Q and X are both matrices  $\in \mathbb{R}^{n \times (p+1)}$ . Hint: this is purely a linear algebra exercise and it's a one-liner.

(h) [easy] Prove that the least squares projection  $\boldsymbol{H} = \boldsymbol{X} \left( \boldsymbol{X}^T \boldsymbol{X} \right)^{-1} \boldsymbol{X}^T = \boldsymbol{Q} \boldsymbol{Q}^{\top}$ . Justify each step.

(i) [difficult] [MA] This problem is independent of the others. Prove that  $\operatorname{rank}[\boldsymbol{H}] = \operatorname{tr}[\boldsymbol{H}]$ . Hint: you will need to use facts about eigenvalues and the eigendecomposition of projection matrices.

(j) [harder] Prove that an orthogonal projection onto the colsp [Q] is the same as the sum of the projections onto each column of Q.

(k)	[easy] Explain why adding a new column to $\boldsymbol{X}$ results in no change in the SST remain ing the same.
(1)	[harder] Prove that adding a new column to $\boldsymbol{X}$ results in SSR increasing.
m)	[harder] What is overfitting? Use what you learned in this problem to frame you answer.
(n)	[easy] Why are "in-sample" error metrics (e.g. $R^2$ , SSE, $s_e$ ) dishonest? Note: I'n leaving out RMSE as RMSE attempts to be honest by increasing as $p$ increases due to the denominator. I've chosen to use standard error of the residuals as the error metric of choice going forward.

(o) [easy] How can we provide honest error metrics (e.g.  $R^2$ , SSE,  $s_e$ )? It may help to draw a picture of the procedure.

(p) [easy] The procedure in (o) produces highly variable honest error metrics. Can you change the procedure slightly to reduce the variation in the honest error metrics? What is this procedure called and how is it done?

#### Problem 3

These are some questions related to polynomial-derived features and logarithm-derived features in use in OLS regression.

(a) [harder] What was the overarching problem we were trying to solve when we started to introduce polynomial terms into  $\mathcal{H}$ ? What was the mathematical theory that justified this solution? Did this turn out to be a good solution? Why / why not?

(b) [harder] We fit the following model:  $\hat{\boldsymbol{y}} = b_0 + b_1 x + b_2 x^2$ . What is the interpretation of  $b_1$ ? What is the interpretation of  $b_2$ ? Although we didn't yet discuss the "true" interpretation of OLS coefficients, do your best with this.

(c) [difficult] Assuming the model from the previous question, if  $x \in \mathcal{X} = [10.0, 10.1]$ , do you expect to "trust" the estimates  $b_1$  and  $b_2$ ? Why or why not?

(d) [difficult] We fit the following model:  $\hat{\boldsymbol{y}} = b_0 + b_1 x_1 + b_2 \ln{(x_2)}$ . We spoke about in class that  $b_1$  represents loosely the predicted change in response for a proportional movement in  $x_2$ . So e.g. if  $x_2$  increases by 10%, the response is predicted to increase by  $0.1b_2$ . Prove this approximation from first principles.

(e) [easy] When does the approximation from the previous question work? When do you expect the approximation from the previous question not to work?

(f) [harder] We fit the following model:  $\ln(\hat{y}) = b_0 + b_1 x_1 + b_2 \ln(x_2)$ . What is the interpretation of  $b_1$ ? What is the approximate interpretation of  $b_2$ ? Although we didn't yet discuss the "true" interpretation of OLS coefficients, do your best with this.

(g) [easy] Show that the model from the previous question is equal to  $\hat{\boldsymbol{y}} = m_0 m_1^{x_1} x_2^{b_2}$  and interpret  $m_1$ .