

# MATH 343 / 643 Homework #3

Professor Adam Kapelner

Due 11:59PM May 18, 2025

(this document last updated 11:14am on Sunday 27<sup>th</sup> April, 2025)

## 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*. For this homework set, read as much as you can online about the topics we covered.

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

Consider the Poisson linear regression model with one feature, time:

$$Y_1, Y_2, \dots, Y_n \mid t_1, t_2, \dots, t_n \stackrel{\text{ind}}{\sim} \text{Poisson}(e^{\beta_0 + \beta_1 t_i})$$

and consider a Bayesian approach to inference.

- (a) [easy] What is the parameter space for the two parameters of interest?
- (b) [easy] Assume a flat prior  $f(\beta_0, \beta_1) \propto 1$ . Find the kernel of the posterior distribution  $f(\beta_0, \beta_1 \mid y_1, \dots, y_n, t_1, \dots, t_n)$ .
- (c) [easy] Find the log of the kernel of the posterior distribution.
- (d) [easy] Find the kernel of the conditional distribution  $f(\beta_0 \mid y_1, \dots, y_n, t_1, \dots, t_n, \beta_1)$ . Is it a brand name rv?
- (e) [easy] Find the kernel of the conditional distribution  $f(\beta_1 \mid y_1, \dots, y_n, t_1, \dots, t_n, \beta_0)$ . Is it a brand name rv?

- (f) [harder] [MA, not covered on the final] Given your answer in (a), the  $\text{Supp}[\beta_0]$ , provide a proposal distribution for the conditional distribution of  $\beta_0$ :

$$q(\beta_0^* | \beta_{0_{t-1}}, y_1, \dots, y_n, t_1, \dots, t_n, \beta_1, \phi) =$$

- (g) [harder] [MA, not covered on the final] Given your answer in (a), the  $\text{Supp}[\beta_1]$ , provide a proposal distribution for the conditional distribution of  $\beta_1$ :

$$q(\beta_1^* | \beta_{1_{t-1}} y_1, \dots, y_n, t_1, \dots, t_n, \beta_0, \phi) =$$

## Problem 2

This question is about basic causality, structural equation models and their visual representation as directed acyclic graphs (DAGs).

- (a) [easy] We run a OLS to fit  $\hat{y} = b_0 + b_1x$  and find there is a statistically significant rejection of  $H_0 : \beta_1 = 0$ . If this test was decided correctly, what do we call the relationship between  $x$  and  $y$ ? (The answer is one word).
- (b) [easy] If this test was decided incorrectly, what do we call the relationship between  $x$  and  $y$ ? (The answer is two words).
- (c) [easy] Draw an example DAG where  $x$  causes  $y$ .
- (d) [easy] Draw an example DAG where  $x$  is correlated to  $y$  but is not causal.
- (e) [easy] Draw an example DAG that can result in a spurious correlation of  $x$  and  $y$ .
- (f) [easy] Draw an example DAG where  $x$  causes  $y$  but its effect is fully blocked by  $z$ .

- (g) [easy] Draw an example DAG where  $x$  causes  $y$  but its effect is partially blocked by  $z$ .
- (h) [easy] Draw an example DAG that results in a Berkson's paradox between  $x$  and  $y_1$ . Denote the collider variable as  $y_2$ .
- (i) [easy] Draw an example DAG that results in a Simpson's paradox between  $x$  and  $y$ . Denote the confounding variable as  $u$ .
- (j) [easy] In the previous Simpson's paradox DAG, provide an example structural equation for  $y$  and provide an example structural equation for  $x$ .
- (k) [easy] Consider observed covariates  $x_1, x_2, x_3$  and phenomenon  $y$ . Draw a realistic DAG for this setting.

### Problem 3

This question is about causal and correlational interpretations for generalized linear models.

- (a) [easy] We run the following model on the `diamonds` dataset where  $y$  is the price of the diamond

```
> summary(lm(price ~ ., diamonds))
```

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	2184.477	408.197	5.352	8.76e-08	***
carat	11256.978	48.628	231.494	< 2e-16	***
cutGood	579.751	33.592	17.259	< 2e-16	***
cutVery Good	726.783	32.241	22.542	< 2e-16	***
cutPremium	762.144	32.228	23.649	< 2e-16	***
cutIdeal	832.912	33.407	24.932	< 2e-16	***
colorE	-209.118	17.893	-11.687	< 2e-16	***
colorF	-272.854	18.093	-15.081	< 2e-16	***
colorG	-482.039	17.716	-27.209	< 2e-16	***
colorH	-980.267	18.836	-52.043	< 2e-16	***
colorI	-1466.244	21.162	-69.286	< 2e-16	***
colorJ	-2369.398	26.131	-90.674	< 2e-16	***
claritySI2	2702.586	43.818	61.677	< 2e-16	***
claritySI1	3665.472	43.634	84.005	< 2e-16	***
clarityVS2	4267.224	43.853	97.306	< 2e-16	***
clarityVS1	4578.398	44.546	102.779	< 2e-16	***
clarityVVS2	4950.814	45.855	107.967	< 2e-16	***
clarityVVS1	5007.759	47.160	106.187	< 2e-16	***
clarityIF	5345.102	51.024	104.757	< 2e-16	***
depth	-63.806	4.535	-14.071	< 2e-16	***
table	-26.474	2.912	-9.092	< 2e-16	***
x	-1008.261	32.898	-30.648	< 2e-16	***
y	9.609	19.333	0.497	0.619	
z	-50.119	33.486	-1.497	0.134	

What is the interpretation of the  $b$  for `carat` (the unit of this feature is “carats”)?

(b) [difficult] What is the interpretation of the  $b$  for `cutIdeal` (note: the reference category for `cut` is `Fair`)?

(c) [easy] We run the following model on the `Pima.tr2` dataset where  $y$  is 1 if the subject had diabetes or 0 if not.

```
> summary(glm(type ~ ., MASS::Pima.tr2, family = "binomial"))
```

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-9.773062	1.770386	-5.520	3.38e-08	***
npreg	0.103183	0.064694	1.595	0.11073	
glu	0.032117	0.006787	4.732	2.22e-06	***
bp	-0.004768	0.018541	-0.257	0.79707	
skin	-0.001917	0.022500	-0.085	0.93211	
bmi	0.083624	0.042827	1.953	0.05087	.
ped	1.820410	0.665514	2.735	0.00623	**
age	0.041184	0.022091	1.864	0.06228	.

What is the interpretation of the  $b$  for `age` (the unit of this feature is age)?

- (d) [easy] What is the interpretation of the  $b$  for `glu` (the unit of this feature is mg/dL) if `glu` is known to be causal?

- (e) [easy] We run the following model on the `philippines` household dataset where  $y$  is the number of people living in a household.

```
> summary(MASS::glm.nb(total ~ ., read.csv("philippines_housing.csv")))
```

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	1.447108	0.088204	16.406	< 2e-16 ***
locationDavaoRegion	-0.011108	0.064367	-0.173	0.86298
locationIlocosRegion	0.053589	0.063284	0.847	0.39711
locationMetroManila	0.074016	0.056731	1.305	0.19201
locationVisayas	0.131151	0.050440	2.600	0.00932 **
age	-0.004896	0.001136	-4.309	1.64e-05 ***
roofPredominantly Strong Material	0.043376	0.052705	0.823	0.41051

What is the interpretation of the  $b$  for `age` (the unit of this feature is years)?

- (f) [easy] We run the following Weibull regression model on the `lung` dataset where  $y$  is survival of the patient.

```
> lung = na.omit(survival::lung)
> lung$status = lung$status - 1 #needs to be 0=alive, 1=dead
> summary(survreg(Surv(lung$time, lung$status) ~
  inst + sex + ph.ecog + ph.karno + wt.loss, lung))
```

	Value	Std. Error	z	p
(Intercept)	7.13673	0.74732	9.55	< 2e-16
inst	0.02042	0.00877	2.33	0.0199
sex	0.39717	0.13852	2.87	0.0041
ph.ecog	-0.69588	0.15463	-4.50	6.8e-06
ph.karno	-0.01558	0.00749	-2.08	0.0376
wt.loss	0.00977	0.00525	1.86	0.0626
Log(scale)	-0.36704	0.07272	-5.05	4.5e-07

What is the interpretation of the  $b$  for `wt.loss` (the unit of this feature is lbs) if `wt.loss` is known to be causal?

- (g) [easy] What is the interpretation of the  $b$  for `ph.ecog` (the unit of this feature is mg/dL) if `ph.ecog` is known to be causal?



## Problem 4

This problem is about controlling values of variables to allow for causal inference.

- (a) [easy] Redraw the “master decision tree” of what to do in every situation beginning with the root node of “Can we assume a DAG?”

- (b) [easy] Explain why controlling / manipulating the values of  $x$  allows for causal inference of  $x$  on  $y$ .
  
- (c) [harder] Explain why a typical observational study (i.e. just collecting data and assembling it into  $\mathbb{D}$ ) cannot allow for causal inference of  $x$  on  $y$ .
  
- (d) [easy] Give an example case (different from the one we spoke about in class) where controlling / manipulating the values of  $x$  is impossible.
  
- (e) [easy] Give an example case (different from the one we spoke about in class) where controlling / manipulating the values of  $x$  is unethical.
  
- (f) [easy] Give an example case (different from the one we spoke about in class) where controlling / manipulating the values of  $x$  is impractical / unaffordable.

- (g) [difficult] Assume in the `diamonds` dataset that the variable `cut` was manipulated by the experimenter prior to assessing the price  $y$ . This isn't absurd since raw diamonds can be cut differently but their color and clarity cannot be altered. Using the linear regression output from the previous problem, what is the interpretation of the  $b$  for `cutIdeal`. The reference category for this variable is `Fair`.

### Problem 5

This problem is about randomized controlled trials (RCTs). Let  $n$  denote the number of subjects, let  $\mathbf{w}$  denote the variable of interest which you seek causal inference for its effect. Here we assume  $\mathbf{w}$  is a binary allocation / assignment vector of the specific manipulation  $w_i$  for each subject (thus the experiment has “two arms” which is sometimes called a “treatment-control experiment” or “pill-placebo trial” or an “AB test”). Let  $\mathbf{y}$  denote the measurements of the phenomenon of interest for each subject and let  $\mathbf{x}_1, \dots, \mathbf{x}_p$  denote the  $p$  baseline covariate measurements for each subject.

- (a) [easy] How many possible allocations are there in this experiment?
- (b) [easy] What are the three advantages of randomizing  $\mathbf{w}$ ? We spoke about two main advantages and one minor advantage.

(c) [easy] In Fisher's Randomization test, what is the null hypothesis? Explain what this really means.

(d) [easy] Explain step-by-step how to run Fisher's Randomization test.

Assume now that Let  $\mathbf{Y} = \beta_0 \mathbf{1}_n + \beta_T \mathbf{w} + \boldsymbol{\varepsilon}$  where  $\varepsilon_1, \dots, \varepsilon_n \stackrel{iid}{\sim}$  mean zero and has homoskedastic variance  $\sigma^2$ .

(e) [easy] What this the parameter of interest in causal inference? What is its name?

(f) [easy] Assume we employ OLS to estimate  $\beta_T$ . We proved previously that OLS estimators are unbiased for any error distribution with mean zero. Find the  $\mathbb{MSE}[B_T]$ .

(g) [easy] Prove that the optimal  $\mathbf{w}$  has equal allocation to each arm.

(h) [easy] Explain how to run an experiment using the *completely randomized design*.

Assume now that Let  $\mathbf{Y} = \beta_0 \mathbf{1}_n + \beta_T \mathbf{w} + \beta_1 \mathbf{x}_{\cdot 1} + \dots + \beta_p \mathbf{x}_{\cdot p} + \boldsymbol{\mathcal{E}}$  where  $\mathcal{E}_1, \dots, \mathcal{E}_n \stackrel{iid}{\sim}$  mean zero and have homoskedastic variance  $\sigma^2$ .

(i) [difficult] Prove that  $B_T$  is unbiased over the distribution of  $\boldsymbol{\mathcal{E}}$  and  $\mathbf{W}$ .

(j) [easy] What is the purpose using a *restricted design*? That is, using a set of allocations that is a subset of the full set of the completely randomized design.

- (k) [harder] Explain how to run an experiment using Fisher's *blocking design* where you block on  $\mathbf{x}_{.1}$ , a factor with three levels and  $\mathbf{x}_{.2}$ , a factor with two levels.

- (l) [easy] What are the two main disadvantages to using Fisher's *blocking design*?

- (m) [easy] Explain how to run an experiment using Student's *rerandomization design* where you let the imbalance metric be

$$\sum_{j=1}^p \frac{|\bar{x}_{jT} - \bar{x}_{jC}|}{s_{x_{jT}}^2/(n/2) + s_{x_{jC}}^2/(n/2)}$$

(n) [easy] Explain how to run an experiment using the *pairwise matching design*.

(o) [easy] Does the pairwise matching design provide better imbalance on the observed covariates than the rerandomization design? Y/N

### Problem 6

This question is about hazard rates and Cox proportion hazard models. This will only be required if we have the two lectures on these models.

(a) [easy] What is the definition of the hazard rate  $h(t)$ ?

(b) [easy] If  $X \sim U(0, 1)$ , derive the hazard rate  $h(t)$ .

(c) [easy] Give an example of a real-world phenomenon  $T$  whose  $h(t)$  is a bathtub shape.

(d) [easy] Prove that  $S(t) = e^{-\int_0^t h(u)du}$ .

(e) [difficult] Explain why the assumption that  $h(t) = h_0(t)e^{\beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}$  is called the “proportional hazard model”.

(f) [easy] Under the proportional hazard model, find the likelihood  $\mathcal{L}(\boldsymbol{\beta}, h_0; \mathbf{X}, \mathbf{y})$ .

(g) [easy] Now let  $h_i := h_0(y_i)$  and  $H_i := \int_0^{y_i} h_0(u) du$ . Find  $\mathcal{L}(\boldsymbol{\beta}, h_1, \dots, h_n, H_1, \dots, H_n; \mathbf{X}, \mathbf{y})$ .

(h) [easy] Now assume (1) all  $y_i$ 's are uniquely-valued and (2)  $H_i \approx h_1 + \dots + h_i$  and find  $\hat{h}_i^{MLE}$ .



(i) [easy] [MA] Find  $\hat{\beta}^{MLE}$ .

(j) [harder] We now run the following Cox proportional hazard model on the `lung` dataset where  $y$  is survival of the patient.

```
> mod = coxph(surv_obj ~ inst + sex + ph.ecog + ph.karno + wt.loss, lung)
> summary(mod)
```

	coef	exp(coef)	se(coef)	z	Pr(> z )	
inst	-0.030042	0.970404	0.012931	-2.323	0.02016	*
sex	-0.571959	0.564419	0.198865	-2.876	0.00403	**
ph.ecog	0.993224	2.699926	0.232115	4.279	1.88e-05	***
ph.karno	0.021492	1.021725	0.011222	1.915	0.05547	.
wt.loss	-0.014800	0.985309	0.007664	-1.931	0.05348	.

What is the interpretation of the  $b$  for `wt.loss` (the unit of this feature is lbs)?