

Math 343 / 643 Spring 2025

Final Examination

Professor Adam Kapelner

May 20, 2025

Full Name _____

Code of Academic Integrity

Since the college is an academic community, its fundamental purpose is the pursuit of knowledge. Essential to the success of this educational mission is a commitment to the principles of academic integrity. Every member of the college community is responsible for upholding the highest standards of honesty at all times. Students, as members of the community, are also responsible for adhering to the principles and spirit of the following Code of Academic Integrity.

Activities that have the effect or intention of interfering with education, pursuit of knowledge, or fair evaluation of a student's performance are prohibited. Examples of such activities include but are not limited to the following definitions:

Cheating Using or attempting to use unauthorized assistance, material, or study aids in examinations or other academic work or preventing, or attempting to prevent, another from using authorized assistance, material, or study aids. Example: using an unauthorized cheat sheet in a quiz or exam, altering a graded exam and resubmitting it for a better grade, etc.

I acknowledge and agree to uphold this Code of Academic Integrity.

signature

date

Instructions

This exam is 120 minutes (variable time per question) and closed-book. You are allowed **three** pages (front and back) of “cheat sheets”, blank scrap paper (provided by the proctor) and a graphing calculator (which is not your smartphone). Please read the questions carefully. Within each problem, I recommend considering the questions that are easy first and then circling back to evaluate the harder ones. No food is allowed, only drinks.

Problem 1 You are trying to model venture capital returns in early-stage startups. Many returns are zero, but of those that are nonzero, they follow a long positive tail. E.g., if you invested in Uber at the seed round, \$1000 would've blossomed to \$5,000,000 at the time of IPO for a 5,000x return. Here, we will model the multiple on investment at time of sale (i.e., this example would have $y = 5000$) and we assume the data follows a “hurdle model”. Why a hurdle? Because a large proportion of startups will fail, returning $y = 0$.

If the investment doesn't fail, we assume its multiple follows a Lomax distribution (which is essentially a ParetoI shifted to the left to begin support at zero) defined below:

$$Y \sim \text{Lomax}(\theta_1, \theta_2) := \frac{\theta_2}{\theta_1} \left(1 + \frac{y}{\theta_1}\right)^{-(\theta_2+1)} \mathbf{1}_{y>0}.$$

This rv model has two parameters, $\theta_1 \in (0, \infty)$ which controls the mean and $\theta_2 \in (0, \infty)$ which scales the mean. Below are the mean and variance.

$$\mathbb{E}[Y] = \begin{cases} \frac{\theta_1}{\theta_2 - 1} & \text{if } \theta_2 > 1 \\ \text{undefined} & \text{otherwise} \end{cases}, \quad \text{Var}[Y] = \begin{cases} \frac{\theta_1^2 \alpha}{(\alpha - 1)^2 (\alpha - 2)} & \alpha > 2 \\ \infty & 1 < \alpha \leq 2 \\ \text{undefined} & \text{otherwise} \end{cases}$$

Thus, the hurdle model for multiples on initial investment is:

$$Y_i \stackrel{\text{ind}}{\sim} \begin{cases} 0 & \text{w.p. } \theta_3 \\ \text{Lomax}(\theta_{1,i}, \theta_2) & \text{w.p. } 1 - \theta_3 \end{cases}$$

As we see above, we will assume θ_1 varies with features of the company at time of investment (as it is indexed by i) but θ_2 and θ_3 do not. A more complex model can explore covariate dependencies in those other two parameters. (That can be a nice masters thesis in finance).

What is our data? We have n observations $\langle \mathbf{x}_1, y_1 \rangle, \dots, \langle \mathbf{x}_n, y_n \rangle$ pairs stacked like in 342 as \mathbf{X}, \mathbf{y} where each subject i has p measurements in row vector $\mathbf{x}_i \in \mathbb{R}^p$ and multiple $y_i \in [0, \infty)$. We wish to use a GLM to specify $\theta_{1,i}$ with the p covariate features so we can fit “linear” coefficients $\boldsymbol{\beta} := [\beta_0 \ \beta_1 \ \dots \ \beta_p]^\top$. We do so with $\theta_{1,i} = \phi(\mathbf{x}_i \mathbf{b}) = e^{\mathbf{x}_i \mathbf{b}}$ just like we saw with Poisson, Negative Binomial and Weibull regression. Also, for convenience, let

$$n_0 := \sum_{i=1}^n \mathbf{1}_{y_i=0} \quad \text{and} \quad n_+ := \sum_{i=1}^n \mathbf{1}_{y_i>0},$$

i.e., the # of startups that failed and the # of startups that did not fail respectively.

(a) [3 pt / 3 pts] How many scalar parameters can we draw inference for in this model?

$p + 1$ for the β_j 's, one for the θ_2 and one for θ_3 yields $p + 3$ total.

- (b) [5 pt / 8 pts] Write the likelihood function from the definition of the hurdle model and then show it can be simplified to the expression at the bottom. Show your work.

$$\begin{aligned}
\mathcal{L}(\boldsymbol{\beta}, \theta_2, \theta_3; \mathbf{X}, \mathbf{y}) &= \prod_{i=1}^n \theta_3^{\mathbb{1}_{y_i=0}} \left((1 - \theta_3) \frac{\theta_2}{e^{\mathbf{x}_i \boldsymbol{\beta}}} \left(1 + \frac{y_i}{\theta_{1,i}} \right)^{-(\theta_2+1)} \right)^{\mathbb{1}_{y_i>0}} \\
&= \prod_{i=1}^n \theta_3^{\mathbb{1}_{y_i=0}} \prod_{i=1}^n (1 - \theta_3)^{\mathbb{1}_{y_i>0}} \prod_{i=1}^n \theta_2^{\mathbb{1}_{y_i>0}} \prod_{i=1}^n \left(\frac{1}{e^{\mathbf{x}_i \boldsymbol{\beta}}} \right)^{\mathbb{1}_{y_i>0}} \prod_{i=1}^n \left(1 + \frac{y_i}{\theta_{1,i}} \right)^{\mathbb{1}_{y_i>0}} \\
&= \theta_3^{n_0} (1 - \theta_3)^{n_+} \theta_2^{n_+} e^{-\left(\sum_{i:y_i>0} \mathbf{x}_i \right) \boldsymbol{\beta}} \left(\prod_{i:y_i>0} 1 + y_i e^{-\mathbf{x}_i \boldsymbol{\beta}} \right)^{-(\theta_2+1)}
\end{aligned}$$

Assume for the rest of the problem that the prior is Laplace i.e. $f(\boldsymbol{\beta}, \theta_2, \theta_3) \propto 1$.

- (c) [4 pt / 12 pts] Find the Gibbs step for θ_3 as a brand-name distribution.

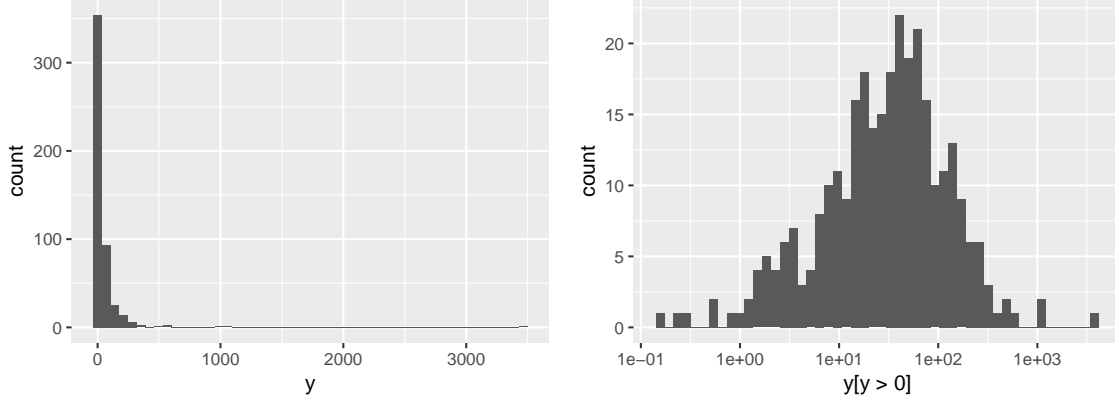
$$f(\theta_3 | \boldsymbol{\beta}, \theta_2, \mathbf{X}, \mathbf{y}) \propto \theta_3^{n_0} (1 - \theta_3)^{n_+} \propto \text{Beta}(n_0 + 1, n_+ + 1)$$

- (d) [2 pt / 14 pts] All the conditional distributions for $\theta_2, \beta_0, \dots, \beta_p$ (given everything else) are not proportional to anyone known distribution. What are the names of the two methods we studied in class to allow for efficient, practical computational Bayesian inference in this scenario?

Metropolis-Hastings sampling, [No U-Turn Sampling within] Hamiltonian MCMC

We now turn to the samples and the features. We have $n = 500$ samples of previous companies with their multiples at exit, \mathbf{y} . We collected $p = 3$ features on each startup at the time of their first funding round: $x_1 := \#$ of cofounders $\in \{1, 2, 3, \dots\}$, $x_2 := 1$ if the startup is the tech space otherwise 0 and $x_3 :=$ seed funding amount (in millions of USD) so it's a positive real number. Thus we have glm parameters $\beta_0, \beta_1, \beta_2, \beta_3$.

Here is a histogram of the raw data on the left and the subset of the raw data where $y > 0$ on the right (i.e. without the zeroes) on a log scale. The $\max(\mathbf{y}) \approx 3500$.

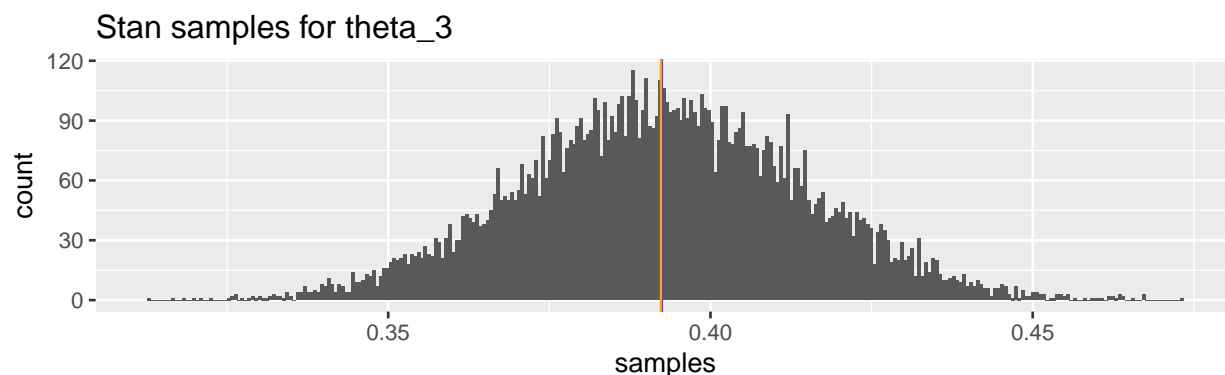


We now use `stan` to do the computational inference.

- (e) [3 pt / 17 pts] To do so, we need to derive the target objective which is the log kernel of the posterior. Find it below:

$$\begin{aligned}
 & \ln(k(\boldsymbol{\beta}, \theta_3, \theta_2 \mid \mathbf{X}, \mathbf{y})) \\
 = & \ln \left(\theta_3^{n_0} (1 - \theta_3)^{n_+} \theta_2^{n_+} e^{-\left(\sum_{i: y_i > 0} \mathbf{x}_i \right) \boldsymbol{\beta}} \left(\prod_{i: y_i > 0} 1 + y_i e^{-\mathbf{x}_i \boldsymbol{\beta}} \right)^{-(\theta_2 + 1)} \right) \\
 = & n_0 \ln(\theta_3) + n_+ \ln(1 - \theta_3) + n_+ \ln(\theta_2) - \left(\sum_{i: y_i > 0} \mathbf{x}_i \right) \boldsymbol{\beta} - (\theta_2 + 1) \sum_{i: y_i > 0} \ln(1 + y_i e^{-\mathbf{x}_i \boldsymbol{\beta}})
 \end{aligned}$$

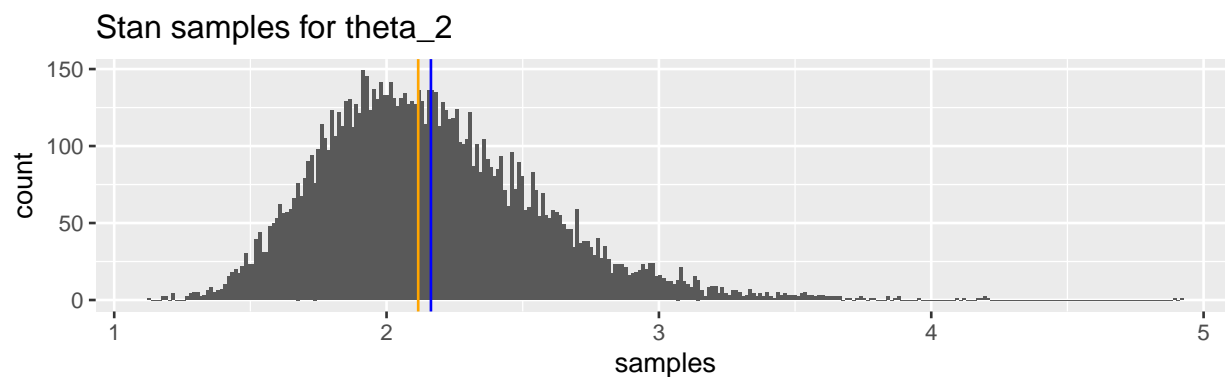
We now do 5,000 samples in `stan`. Below is a histogram of the samples for θ_3 produced by our visualize chains functions we used in class:



- (f) [3 pt / 20 pts] Find a 95% credible region for θ_3 .

$$CR_{\theta_3, 95\%} = [0.35, 0.44]$$

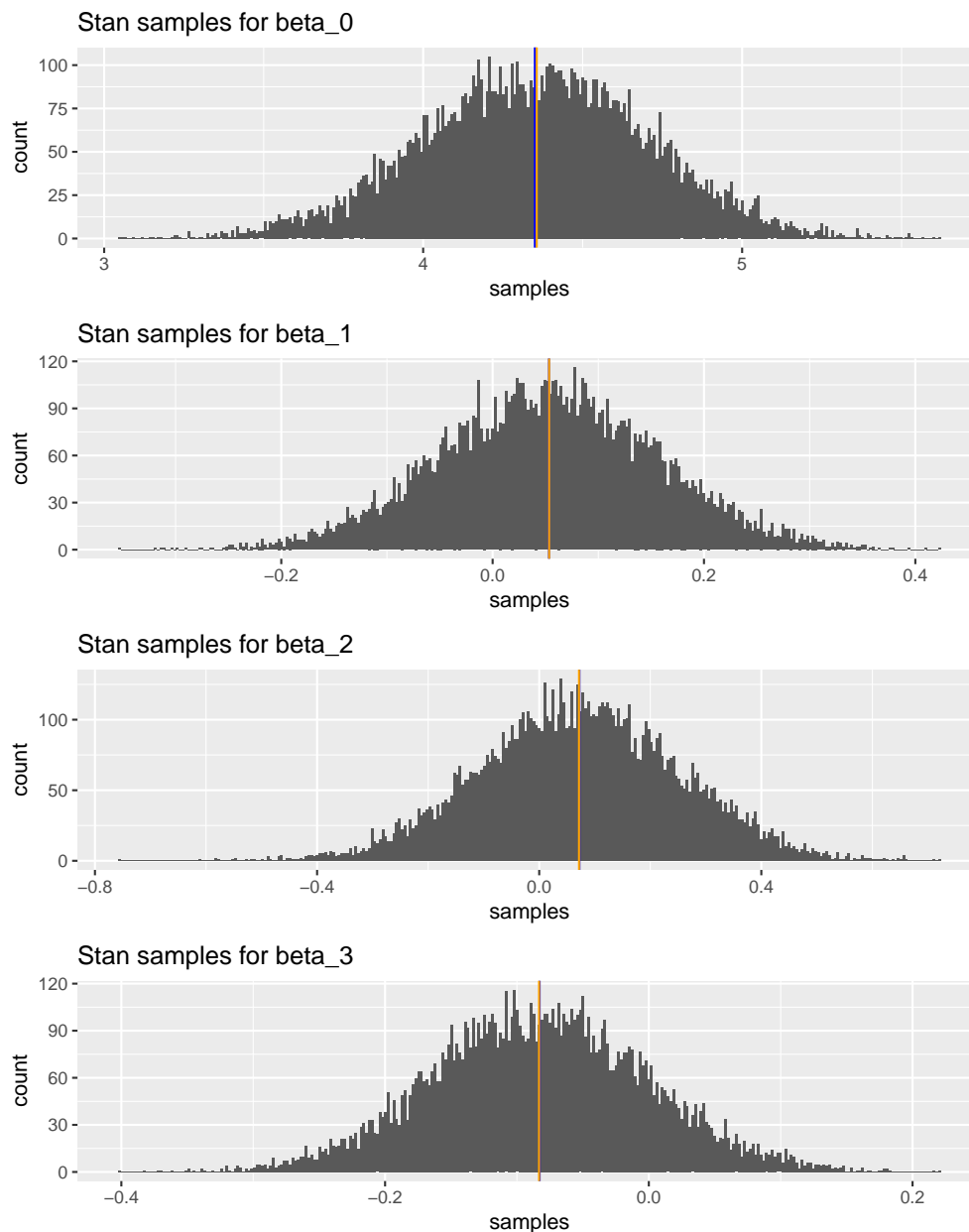
Below is a histogram of the samples for θ_2 produced by our visualize chains functions we used in class:



- (g) [4 pt / 24 pts] The Lomax is a very interesting distribution. If $\theta_2 < 2$, the distribution has infinite or undefined variance (see problem header). Given the data we've seen, what is the probability our DGP has infinite or undefined variance assuming this model?

$$\mathbb{P}(\theta_2 < 2 \mid \mathbf{X}, \mathbf{y}) \approx 40\%.$$

Below are histograms of the samples for all β_j 's produced by our visualize chains functions we used in class:



- (h) [5 pt / 29 pts] Make a decision about the omnibus test (i.e., the null hypothesis is that none of the three features matter in the linear model). Justify your decision.

Since (a) $H_0 : \beta_1 = 0$, $H_0 : \beta_2 = 0$ and $H_0 : \beta_3 = 0$ will be retained at any reasonable δ or (b) CR's for $\beta_1, \beta_2, \beta_3$ will all include zero at any reasonable α , we can assume that we must also retain the omnibus null hypothesis that $H_0 : \beta_1 = \beta_2 = \beta_3 = 0$.

Regardless of what you answered in (h), you now want to examine the following company where $\mathbf{x}_\star = [2 \ 1 \ 5]$ as they have 2 cofounders, it's a tech company and they raised \$5M in their seed round.

- (i) [3 pt / 32 pts] Estimate the probability this company will fail (i.e., return zero).

This is our best guess of θ_3 , the hurdle probability which is unaffected by \mathbf{x} . See top illustration on previous page for the MCMC samples. The best guess of θ_3 given by the vertical line which is both the average value and median value is $\approx 39\%$.

- (j) [5 pt / 37 pts] Assuming this company does not fail, what is the prediction of their return multiple, y_\star rounded to two decimals?

This is the expectation of the Lomax whose formula is given in the problem header. We need the link function, all the β_j estimates and θ_2 :

$$\begin{aligned} y_\star = \mathbb{E}[Y \mid \mathbf{x}_\star] &= \frac{\theta_{1,\star}}{\theta_2 - 1} = \frac{e^{\mathbf{x}_\star \mathbf{b}}}{\theta_2 - 1} = \frac{e^{b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3}}{\theta_2 - 1} \\ &= \frac{\exp(4.3 + 0.05(2) + 0.1(1) + -0.1(5))}{2.1 - 1} = 49.63 \end{aligned}$$

Problem 2 Consider a design matrix \mathbf{X} of size $n \times (p + 1)$ whose first column is $\mathbf{1}_n$ and whose rows are real measurements corresponding to response values in the vector $\mathbf{y} \in \mathbb{R}^n$. We have reason to believe most of these features are uninformative to the response. We use the following algorithms to generate linear coefficients:

$$\begin{aligned} \mathcal{A}_1 : \mathbf{b}_{\mathcal{A}_1} &= \arg \min_{\mathbf{w} \in \mathbb{R}^{p+1}} \{(\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})\} \\ \mathcal{A}_2 : \mathbf{b}_{\mathcal{A}_2} &= \arg \min_{\mathbf{w} \in \mathbb{R}^{p+1}} \left\{ (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \sum_{j=1}^p w_j^2 \right\} \quad \text{where } \lambda > 0 \\ \mathcal{A}_3 : \mathbf{b}_{\mathcal{A}_3} &= \arg \min_{\mathbf{w} \in \mathbb{R}^{p+1}} \left\{ (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) + \gamma \sum_{j=1}^p |w_j| \right\} \quad \text{where } \gamma > 0 \end{aligned}$$

- (a) [3 pt / 40 pts] What are the names of these three algorithms?

\mathcal{A}_1 : OLS
 \mathcal{A}_2 : Ridge regression
 \mathcal{A}_3 : Lasso regression

- (b) [2 pt / 42 pts] For algorithms \mathcal{A}_2 and \mathcal{A}_3 , what is a recommended prestep to do before running the algorithms?

Standardizing the values in columns 2, ..., $p + 1$ to have average zero and standard deviation one.

- (c) [5 pt / 47 pts] Circle all the following statements that are true.

- TRUE If $\lambda = 0$, $\|\mathbf{b}_{\mathcal{A}_1}\| = \|\mathbf{b}_{\mathcal{A}_2}\|$
- $\exists \lambda > 0$, $\|\mathbf{b}_{\mathcal{A}_1}\| = \|\mathbf{b}_{\mathcal{A}_2}\|$
- $\exists \gamma > 0$, $\|\mathbf{b}_{\mathcal{A}_1}\| = \|\mathbf{b}_{\mathcal{A}_3}\|$
- TRUE $\exists \lambda > 0, \gamma > 0$, $\|\mathbf{b}_{\mathcal{A}_2}\| = \|\mathbf{b}_{\mathcal{A}_3}\|$
- TRUE $\lim_{\lambda \rightarrow \infty} \|\mathbf{b}_{\mathcal{A}_2}\| = 0$

- (d) [3 pt / 50 pts] Below are three columns each a sample of entries in \mathbf{b} for some of the p features. Match the most likely algorithms (1, 2, 3) to each of the columns by filling in the lines below the columns with a permutation of $\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3$.

0.000000	-0.3600379	-0.2044754
0.000000	-0.0497710	-0.1043836
0.000000	2.2921732	3.7036262
0.000000	-0.0756272	-0.4866566
0.000000	-5.7274435	-7.4075626
2.866342	3.1096124	3.1156581
0.000000	0.3457908	0.2125626
0.000000	-3.1055237	-3.5124361
0.000000	3.7399084	5.4601868
0.000000	-2.6364656	-4.6368675
-1.352798	-3.6628336	-4.6743968
0.1619854	-0.1791046	-0.1105668
-3.440234	-4.5821048	-4.0857243
0.0000000	0.2331904	0.6371670
0.0000000	-0.1589228	-0.4134656
0.0000000	0.5340268	0.6212814
0.0000000	-0.2324948	-0.4902040
0.0000000	0.1720569	0.3553368
0.0000000	-1.7887053	-2.1417532
0.0000000	0.2537185	0.0777026

\mathcal{A}_3

\mathcal{A}_2

\mathcal{A}_1

Problem 3 For each of the following, you will be provided with a DAG that includes observed metrics x_1, x_2 and an observed response y . These DAGs should be considered complete and self-contained, i.e. there are no other variables in the system and thus no noise. You will then be shown a series of predictive models with with OLS on a simple random sample of data points $\langle x_{1,1}, x_{2,1}, y_1 \rangle, \dots, \langle x_{1,n}, x_{2,n}, y_n \rangle$.

For each problem, circle all coefficient estimates that are *unbiased* estimates of the *causal* effect of x_1 on y (and/or x_2 on y). Note: a causal estimand of zero is still a real estimand. Here is an example done for you:



$$\hat{y} = b_0 + \textcircled{b_1} x_1$$

(a) [4 pt / 54 pts]

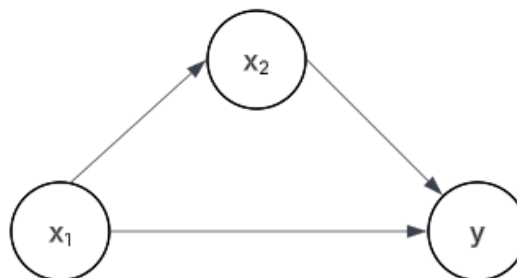


$$\hat{y} = b_0 + \textcolor{red}{b_1} x_1$$

$$\hat{y} = b_0 + \textcolor{red}{b_2} x_2$$

$$\hat{y} = b_0 + b_1 x_1 + \textcolor{red}{b_2} x_2$$

(b) [4 pt / 58 pts]



$$\hat{y} = b_0 + \textcolor{red}{b_1} x_1$$

$$\hat{y} = b_0 + b_2 x_2$$

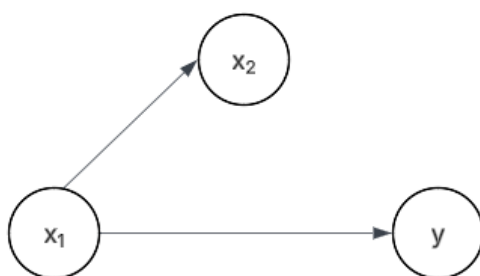
$$\hat{y} = b_0 + b_1 x_1 + \textcolor{red}{b_2} x_2$$

(c) [1 pt / 59 pts]



$$\hat{y} = b_0 + b_1 x_1$$

(d) [4 pt / 63 pts]

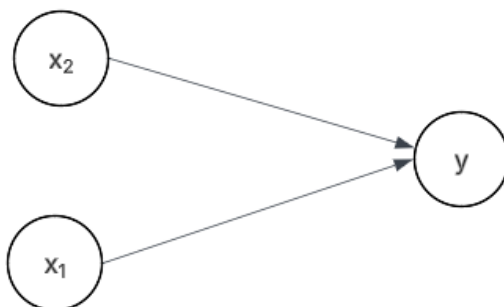


$$\hat{y} = b_0 + b_1 x_1$$

$$\hat{y} = b_0 + b_2 x_2$$

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2$$

(e) [4 pt / 67 pts]



$$\hat{y} = b_0 + b_1 x_1$$

$$\hat{y} = b_0 + b_2 x_2$$

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2$$

Problem 4 The following questions are on experimental design. Let each of the rows be a subject in an experiment the $\mathbf{x}_{\cdot 1}$ column be the only covariate observed before the experiment begins. For each question, create two *different* allocations $\mathbf{w}_1, \mathbf{w}_2$ from the _____ design for this experiment.

- (a) [2 pt / 69 pts] CRD Any binary vectors will be correct here.

$\mathbf{x}_{\cdot 1}$	\mathbf{w}_1	\mathbf{w}_2
0.03		
-0.74		
0.19		
-1.80		
1.47		
0.15		
2.17		
0.48		

- (b) [4 pt / 73 pts] BCRD As long as there are four zeroes and four ones in both columns, it will be correct.

$\mathbf{x}_{\cdot 1}$	\mathbf{w}_1	\mathbf{w}_2
0.03		
-0.74		
0.19		
-1.80		
1.47		
0.15		
2.17		
0.48		

- (c) [6 pt / 79 pts] blocking with $B = 2$ We have a continuous covariate so we sort it to obtain block #1 of the smallest four values $-1.80, -0.74, 0.03, 0.15$, and block #2 of the highest four values $0.19, 0.48, 1.47, 2.17$. The first block will be the subjects corresponding to the smallest four values and the second block to be the subjects corresponding to the highest four values. It is correct if there are two zeroes and two ones in each block.

$\mathbf{x}_{\cdot 1}$	\mathbf{w}_1	\mathbf{w}_2
0.03		
-0.74		
0.19		
-1.80		
1.47		
0.15		
2.17		
0.48		

- (d) [6 pt / 85 pts] pairwise matching The sorted values are $-1.80, -0.74, 0.03, 0.15, 0.19, 0.48, 1.47, 2.17$ so the first pair consist of the subjects that have the covariates measurements $-1.80, -0.74$, the second pair consists of the subjects that have the covariates measurements $0.03, 0.15$, etc. As long as there is one zero and one one in each pair, it is correct.

\mathbf{x}_1	\mathbf{w}_1	\mathbf{w}_2
0.03		
-0.74		
0.19		
-1.80		
1.47		
0.15		
2.17		
0.48		

Problem 5 The following questions are on based on glm's.

- (a) [5 pt / 90 pts] Consider the following code and snippet of its output:

```
> pima = na.omit(MASS::Pima.tr2)
> pima$type = ifelse(pima$type == "Yes", 1, 0)
> summary(glm(type ~ ., pima, family = binomial(link = "logit")))
```

	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	

Assume the data was collected using a simple random sample of subjects and that the `type` being equal to “Yes” means the subject has diabetes otherwise not. Provide an interpretation of the estimated coefficient of `glu` which is measured in mg/dL.

When comparing two subjects (A) and (B) which are sampled in the same fashion as the other subjects in this dataset where (A) has a `glu` measurement 1mg/dL larger than (B)'s `glu` measurement but share the same other observed measurements otherwise, then (A) is predicted to have an estimated log odds probability of diabetes 0.032 ± 0.007 higher than (B)'s log odds probability of diabetes assuming the log odds probability of diabetes is linear in the measurements considered herein.

(b) [5 pt / 95 pts] Consider the following code and snippet of its output:

```
> lung = na.omit(survival::lung)
> lung$status = lung$status - 1 #needs to be 0=alive, 1=dead
> surv_obj = Surv(lung$time, lung$status)
> full_mod = survreg(surv_obj ~ . - time - status, lung)
> summary(full_mod)
```

	Value	Std. Error	z	p
(Intercept)	7.17e+00	1.08e+00	6.64	3.2e-11
inst	2.05e-02	8.80e-03	2.32	0.0201
age	-7.54e-03	8.06e-03	-0.94	0.3497
sex	3.91e-01	1.39e-01	2.82	0.0048
ph.ecog	-6.26e-01	1.58e-01	-3.95	7.7e-05
ph.karno	-1.86e-02	7.67e-03	-2.43	0.0151
pat.karno	7.68e-03	5.54e-03	1.39	0.1656
meal.cal	7.32e-06	1.81e-04	0.04	0.9678
wt.loss	1.10e-02	5.34e-03	2.07	0.0387
Log(scale)	-3.76e-01	7.28e-02	-5.17	2.3e-07

Scale= 0.686

Assume the data was collected using a simple random sample of subjects, and the variable `ph.ecog` (which is measured in “score units”) was experimentally manipulated in the subjects while all other variables were observed naturally. The variable `time` is measured in units of days. Provide an interpretation of the estimated coefficient of `ph.ecog`.

If `ph.ecog` is increased by one score unit and all other measurements remain constant, the [log survival of the subject in days will resultingly decrease by an estimated 0.626 ± 0.158] assuming subject survival-in-days is Weibull-distributed with log mean linear in the measurements considered herein.

You can alternatively replace the above bracketed text with “the survival of the subject in days will resultingly decrease by 46.5%” where that value is calculated as $1 - e^{-0.626}$.

Problem 6 This is the same setup as problem 1 from Midterm II. Consider the following full-rank design matrix:

$$\mathbf{X} := [\mathbf{1}_n \mid \mathbf{x}_{\cdot 1} \mid \dots \mid \mathbf{x}_{\cdot p}] = \begin{bmatrix} \mathbf{x}_{1\cdot} \\ \vdots \\ \mathbf{x}_{n\cdot} \end{bmatrix}$$

with column indices $0, 1, \dots, p$ and row indices $1, 2, \dots, n$. Let \mathbf{H} be the orthogonal projection matrix onto the column space of \mathbf{X} and let $\mathbf{I}_n - \mathbf{H}$ be the orthogonal projection matrix onto the column space of \mathbf{X}_\perp . We assume also a continuous (real-valued) response model which is linear in these measurements, i.e. $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$. For the error term, we assume the “core assumption”,

$$\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}_n, \sigma^2 \mathbf{I}_n) \quad \text{where } \sigma^2 > 0.$$

Consider the following estimator for $\boldsymbol{\beta}$: $\mathbf{B} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ and let $\hat{\mathbf{Y}} := \mathbf{X}\mathbf{B}$ and $\mathbf{E} := \mathbf{Y} - \hat{\mathbf{Y}}$.

- (a) [5 pt / 100 pts] Derive the distribution of \mathbf{E} with only what is in the problem header, the fact about multivariate normal distributions from 340 and linear algebra manipulations. Show each step.

$$\begin{aligned} \mathbf{E} &= (\mathbf{I}_n - \mathbf{H})\mathbf{Y} \\ &= (\mathbf{I}_n - \mathbf{H})(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) \\ &= (\mathbf{I}_n - \mathbf{H})\mathbf{X}\boldsymbol{\beta} + (\mathbf{I}_n - \mathbf{H})\boldsymbol{\varepsilon} \\ &= (\mathbf{I}_n - \mathbf{H})\boldsymbol{\varepsilon} \\ &\sim \mathcal{N}_n(\mathbf{0}_n, (\mathbf{I}_n - \mathbf{H})\sigma^2 \mathbf{I}_n (\mathbf{I}_n - \mathbf{H})^\top) \\ &= \mathcal{N}_n(\mathbf{0}_n, \sigma^2(\mathbf{I}_n - \mathbf{H})) \end{aligned}$$