

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?

$$\beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}$$

- (b) [easy] Assume a flat prior $f(\beta_0, \beta_1) \propto 1$. Find the kernel of the posterior distribution $f(\beta_0, \beta_1 \mid y_1, \dots, y_n, t_1, \dots, t_n)$.

$$\prod e^{-e^{\beta_0 + \beta_1 t_i}} e^{(y_i - (\beta_0 + \beta_1 t_i))}$$

~~y_i~~

- (c) [easy] Find the log of the kernel of the posterior distribution.

$$\sum -e^{\beta_0 + \beta_1 t_i} + y_i (\beta_0 + \beta_1 t_i)$$

- (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?

$$\{ (-e^{\beta_0 + \beta_1 t_i} + y_i (\beta_0 + \beta_1 t_i))$$

$$f(\beta_0 \mid \cdot) \propto e$$

If is not 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?

$$\{ (-e^{\beta_0 + \beta_1 t_i} + y_i \beta_1 t_i)$$

$$f(\beta_1 \mid \cdot) \propto e$$

If is not 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 β_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 β_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).

correlation

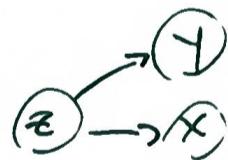
- (b) [easy] If this test was decided incorrectly, what do we call the relationship between x and y ? (The answer is two words).

spurious correlation

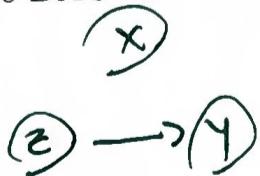
- (c) [easy] Draw an example DAG where x causes y .



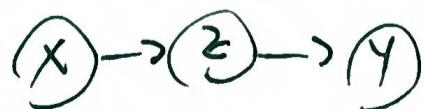
- (d) [easy] Draw an example DAG where x is correlated to y but is not causal.



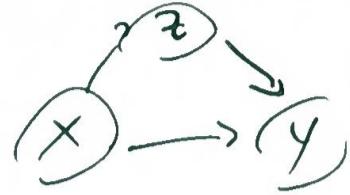
- (e) [easy] Draw an example DAG that can result in a spurious correlation of x and y .



- (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 .

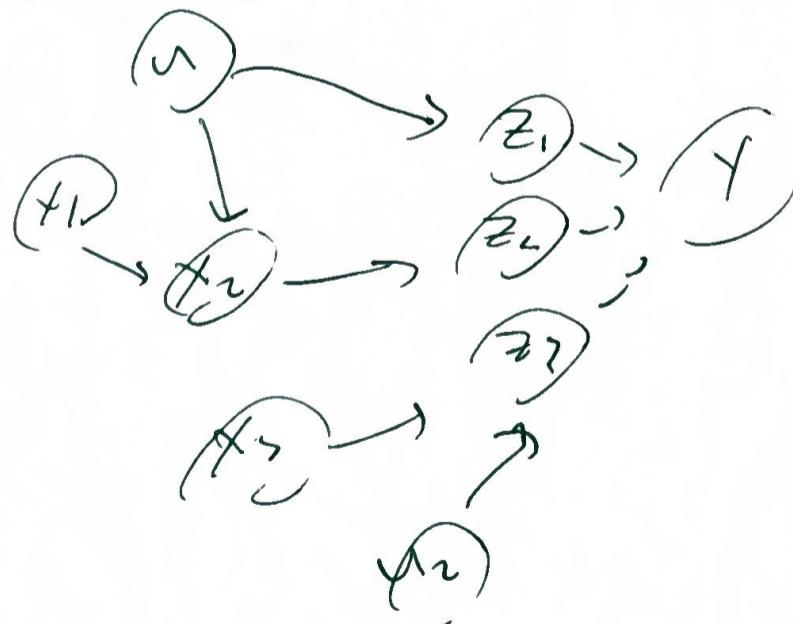


- (j) [easy] In the previous Simpson's paradox DAG, provide an example structural equation for y and provide an example structural equation for x .

$$x = \beta_0 + \beta_1 u + \epsilon$$

$$y = \beta_0 + \beta_1 u + \beta_2 x + \epsilon$$

- (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")?

When comparing two observations (A) and (B) which are sampled in the same fashion as D where (A) has carat 1 carats larger than B but share all other measurement values, then (A) is predicted to have a estimated price $11256.978 + 48.628 \cdot 1$ dollars higher than (B)'s price assuming price is linear in the covariates and the model is stationary

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

When comparing two obs (A) and (B) which are sampled in the same fashion as ID (A) has cut^{1st} 1 at ~~0.01~~ larger than (B) ~~and~~ respective to cutFair but share all other measurement values, then (A) is predicted to have an estimated price 832.91 ± 33.407 dollars higher than (B)'s price assuming price is linear in the P covariants and the model is stationary

- (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)?

When comparing two obs (A) and (B) which are sampled in the same fashion as ID where (A) has ~~been~~ ~~been~~ age 1 year larger than (B) but share all other measurement values. Then (A) is predicted to have an estimated log-odds of having diabetes of $.0411 \pm .022$ larger than (B) assuming log-odds of diabetes is linear in the P covariants and the model is stationary

- (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?

If glucose blood sugar is measured by one mg/dL and all other measurements remain constant, the log odds of getting diabetes will resultingly increase by an estimated $.0321 + .0067$
assuming log odds of getting diabetes is linear in the p covariates and the relationships remain stationary.

- (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)?

If age is measured by one year and all other measurements remain constant, the num. of people living in a household

then compare for obs's (A) and (B) which are sampled in the same fashion as D where (A) has older age than (B) but share all the measurement values from (A). It is predicted to have an estimated decrease of Num. of people living in a household by $e^{-0.048}$ assuming age is multiplicative in the p covariates

- (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?

If wt.loss is inc by one lbs and all other measures remain constant, the log survival years will resultingly decrease by $-0.977 \pm .0052$ assuming survival in Weibull-distribution with log linear mean linear in the p covariates and the relationship is stationary

- (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?

If ph.ecog is inc. by one mg/dL and all other measures remain constant, the log survival years will resultingly decrease by $.695 \pm .154$ assuming survival in Weibull-distribution with log mean linear in the p covariates and the relationship is stationary

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 .

Because it separates x from other
"nodes"

- (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 .

Because there is no manipulation happening

- (d) [easy] Give an example case (different from the one we spoke about in class) where controlling / manipulating the values of x is impossible.

$X = \text{age}$

- (e) [easy] Give an example case (different from the one we spoke about in class) where controlling / manipulating the values of x is unethical.

$X = \text{being homeless}$

- (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.

$X = \text{traveling on private jet}$

- (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.

If cut is measured by one cut respective to ~~the~~ cutIdeal
and all other measurements remain constant, price will
resultingly increase by an
estimated $8.32.92 \pm 37.407$ dollars assuming
price is linear in the p covariate and the
relationship remains stationary

Problem 5

This problem is about randomized controlled trials (RCTs). Let n denote the number of subjects, let w denote the variable of interest which you seek causal inference for its effect. Here we assume 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 y denote the measurements of the phenomenon of interest for each subject and let x_1, \dots, x_p denote the p baseline covariate measurements for each subject.

- (a) [easy] How many possible allocations are there in this experiment? $\binom{n}{n/2}$
(b) [easy] What are the three advantages of randomizing w ? We spoke about two main advantages and one minor advantage.

① On average, over all randomized experiments,
 B_1 is unbiased

② $\bar{y}_1 - \bar{y}_0$ is small, approaches 0 as n gets bigger

③ Reduced bias for inference

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

$$H_0: \gamma_i(w_i=1) = \gamma_i(w_i=0) \quad \forall i$$

for every single subject
the treatment w has no effect

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

- ① run experiment with \vec{w}_{exp} and get \bar{T}
- ② generate \vec{w} from \vec{W}
- ③ compute $b_T = \bar{y}_1 - \bar{y}_0$ for many \vec{w}_i and construct Null dist for B_T
- ④ If $b_{\text{exp}} \notin [Q(\frac{\alpha}{2}), Q(1-\frac{\alpha}{2})]$ \Rightarrow Rehn H_0
else Reject

Assume now that Let $\mathbf{Y} = \beta_0 \mathbf{1}_n + \beta_T \mathbf{w} + \boldsymbol{\varepsilon}$ where $\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_n \stackrel{iid}{\sim}$ mean zero and has homoskedastic variance σ^2 .

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

B_T , the PATE

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

$$\begin{aligned} \text{MSE}[\beta_T] &= 6^2 (\mathbf{X}' \mathbf{X}^{-1}) \\ &= 6^2 \cdot \frac{1}{n} \frac{1}{1-p_1} \begin{bmatrix} 1 & -1 \\ -1 & p_1 \end{bmatrix} \end{aligned}$$

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

$$\begin{aligned}\vec{w}_* &= \arg \min \left\{ \frac{1}{2} (\vec{X}^\top \vec{X})^{-1} \vec{Z} \right\} \\ &= \arg \min \left\{ \frac{1}{n} \left(\frac{1}{1-p_1} + \frac{1}{1-p_2} \right) \right\} \\ &= \arg \max \left\{ p_1 (1-p_1) \right\} \\ &= \left\{ \vec{w}: p_1 = \frac{1}{2} \right\}\end{aligned}$$

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

assign \vec{w} as from a realization of $\vec{w} \sim \mathcal{N} \left(\vec{0}, \frac{1}{n} \mathbb{I}_n \right)$

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

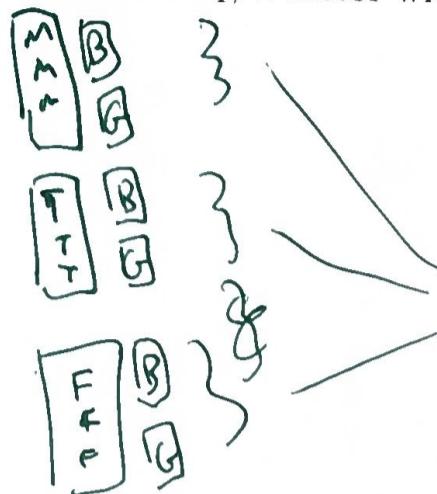
(i) [difficult] Prove that B_T is unbiased over the distribution of $\mathbf{\epsilon}$ and \mathbf{W} .

$$\begin{aligned}B_T &= \frac{2 \vec{w}^\top \vec{y} - \vec{1}^\top \vec{y}}{n/2} = \mathbb{E}_\epsilon \left(\vec{w}^\top (\beta_0 \vec{1}_n + \beta_T \vec{w} + \beta_1 \vec{u} + \vec{\epsilon}) - \frac{2}{n} \vec{1}_n^\top (\beta_0 \vec{1}_n + \beta_T \vec{w} + \beta_1 \vec{u} + \vec{\epsilon}) \right) \\ &= \frac{1}{n} (\beta_0 \vec{1}_n^\top \vec{w} + \beta_T \vec{w}^\top \vec{w} + \beta_1 \vec{u}^\top \vec{w} + \vec{w}^\top \vec{\epsilon}) - \frac{2}{n} (\beta_0 \vec{1}_n^\top \vec{1}_n + \beta_T \vec{w}^\top \vec{1}_n + \beta_1 \vec{u}^\top \vec{1}_n + \vec{1}_n^\top \vec{\epsilon}) \\ &= (2\beta_0 - 2\beta_1) - \frac{4\beta_T}{n} \vec{w}^\top \vec{w} + 2\beta_1 + \frac{4\beta_u}{n} \vec{u}^\top \vec{w} - \frac{2\beta_u}{n} \vec{1}_n^\top \vec{w} - \beta_1 - \frac{2}{n} \vec{1}_n^\top \vec{\epsilon} \\ &= (2\beta_0 - 2\beta_1) - \frac{4\beta_T}{n} \vec{w}^\top \vec{w} + 2\beta_1 + \frac{2\beta_u}{n} (2\vec{u}^\top \vec{w} - \vec{w}^\top \vec{1}_n) \\ \mathbb{E}_\epsilon [B_T] &= \beta_0 + \frac{4\beta_u}{n} \vec{u}^\top \vec{w} - \frac{2\beta_u}{n} \vec{1}_n^\top \vec{w} = \beta_0 + \frac{2\beta_u}{n} (2\vec{u}^\top \vec{w} - \vec{w}^\top \vec{1}_n) \\ \mathbb{E}_w [\mathbb{E}_\epsilon [B_T]] &= \beta_0 + \frac{2\beta_u}{n} (2\vec{u}^\top \vec{w} + \vec{w}^\top \vec{1}_n) \\ &= \beta_0 + 2\frac{\beta_u}{n} \left(2\vec{u}^\top \left(\frac{1}{2} \vec{1}_n \right) - \vec{u}^\top \vec{1}_n \right) = \beta_0\end{aligned}$$

(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.

To reduce bias stemming from
the covariates

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



$$\text{Let } x_{.1} = \{M, T, P\}$$

$$\text{Let } x_{.2} = \{B, G\}$$

$$\text{set } \bar{w} \text{ s.t. } n_1 = n_C$$



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

- ① After a few ~~few~~ factors blocked, # blocks increases exp.
- ② \bar{w} is distributed less randomly

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

- ① define distance function $d(X, \bar{w}) = \sum_{j=1}^p \frac{|\bar{x}_{jT} - \bar{x}_{jC}|}{s_{x_{jT}}^2 / (n/2) + s_{x_{jC}}^2 / (n/2)}$
- ② if $d(X, \bar{w}) > d_{th}$ (hyperparameter), then draw \bar{w} again

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

① normalize X s.t. all covariates have $\text{avg} \approx 0, \text{sd} \approx 1$

② define distance function $d(\vec{x}_i, \vec{x}_j)$

③ Calc D the norm upper-triangular matrix
of all $d(\vec{x}_i, \vec{x}_j)$

④ Compute pairs $\{(i(1), j(1)), \dots, \{(i(n), j(n))\}$
s.t. $d(x_{i(1)}, x_{j(1)}), d(x_{i(2)}, x_{j(2)}) \dots$
 $, d(x_{i(n)}, x_{j(n)})$ is minimum

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