1. The objective of this problem is to demonstrate that *the moments may not always determine a probability distribution.* Consider a random variable $X$ having the standard lognormal distribution with density given by

$$f(x) = \begin{cases} \frac{1}{x\sqrt{2\pi}} \exp\{-\frac{1}{2}(\log x)^2\} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

(a) *(3 points)* Show that the $k$-th raw moment of $X$, corresponding to the above density function $f(x)$ is given by

$$E(X^k) = \exp\{\frac{1}{2}k^2\}, \quad \text{for } k = 1, 2 \ldots$$

(b) *(3 points)* Consider the function given by

$$g(y) = \begin{cases} f(y)(1 + \sin(2\pi \log(y))) & \text{if } y > 0 \\ 0 & \text{otherwise,} \end{cases}$$

where $f(y)$ is as defined above. Show that $g(y)$ is a probability density function.

(c) *(4 points)* Let $Y$ be a random variable with density function $g(y)$ as defined in part (b). Show that

$$E(Y^k) = E(X^k), \quad \text{for } k = 1, 2 \ldots$$

Problem 2

(a) $M_X(t) = E(e^{tx})$

$E(x^k) = \int_0^\infty x^k \frac{1}{x\sqrt{2\pi}} \exp(-\frac{1}{2}(\log x)^2) dx$

~~$= \int_0^\infty e^{tx} \frac{1}{x\sqrt{2\pi}} \exp\{-\frac{1}{2}(\log x)^2\} dx$~~

~~$= \frac{1}{\sqrt{2\pi}} \int_0^\infty \frac{1}{x} \exp\{tx - \frac{1}{2}(\log x)^2\} dx$~~

$= \int_0^\infty \frac{1}{\sqrt{2\pi}} x^{k-1} \exp(-\frac{1}{2}(\log x)^2) dx$ ✗

$= \int_0^\infty \frac{1}{\sqrt{2\pi}} x^{k+1} \exp(-\frac{1}{2}(\log x)^2 + $ ✗

$E(x^k) = \frac{\partial^k}{\partial t^k} M_X(t)|_{t=0}$

$= \frac{\partial^{k-1}}{\partial t^{k-1}} \frac{1}{\sqrt{2\pi}} \int_0^\infty \frac{x}{x} \exp\{tx - \frac{1}{2}(\log x)^2\} dx |_{t=0}$

$= \frac{1}{\sqrt{2\pi}} \int_0^\infty x^{k-1} \exp\{tx - \frac{1}{2}(\log x^2)\} dx|_{t=0}$

$\boxed{=} \exp(\frac{1}{2}x^2)$

*[left margin:]* $\frac{\partial^k}{\partial t^k} M_X(t)$ is not finite

(b) $\int_0^\infty g(y) dy$

$= \int_0^\infty \frac{1}{y\sqrt{2\pi}} \exp\{-\frac{1}{2}(\log y)^2\} (1 + \sin(2\pi \log(y))) dy$

$= \int_0^\infty \frac{1}{y\sqrt{2\pi}} \exp\{-\frac{1}{2}(\log y)^2\} dy + \boxed{\int_0^\infty \frac{1}{y\sqrt{2\pi}} \exp\{-\frac{1}{2}(\log y)^2\} \sin(2\pi \log(y)) dy}$
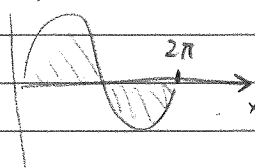
$= 1 \boxed{(+0)}$

$= 1$

Thus $g(y)$ is a probability function, as it sums to 1 and is 0 everywhere else outside its support.

(this since $f(y) = 0$ when $y < 0$).

(Thus $g(y) = 0$ when $y < 0$)

Why $g(y) \geq 0$?

$\boxed{-1}$

*[right side, sin(x) sketch]*

$\sin(x)$ with shaded region, marked $2\pi$, axis $x$

This is because recall $\sin(\cdot)$ is symmetric to 0, thus the sum (integral) due to the symmetry, equals to 0.

Thus $= E(\sin(2\pi \log(y)))$ where $y \sim$ lognormal

$= 0$

*[bottom section]*

$= \frac{1}{\sqrt{2\pi}} \int_0^\infty x^{k-1} \exp\{-\frac{1}{2}(\log(x))^2\} dx$     at $t = 0$

$= \frac{1}{\sqrt{2\pi}} \int_0^\infty e^{uk} \exp\{-\frac{1}{2}u^2\} du$     let $u = \log(x)$  $du = dx \frac{1}{x}$

$= E(e^{uk})$ where $u \sim N(0,1)$

$= e^{\frac{1}{2}k^2}$ ⟵ MGF of $N(0,1)$

Problem 1 continuous

(c). Showing $E(y^k) = E(x^k)$

is equivalent to show $E(y^k) = e^{\frac{1}{2}k^2}$

$M_y(t) = E(e^{ty})$

$= \int_0^\infty e^{ty} f(y)(1+\sin(2\pi \log(y))) \, dy$

$= \int_0^\infty e^{ty} f(y) \, dy + \int_0^\infty e^{ty} \underbrace{\sin(2\pi \log(y))}_{f(y)} \, dy$

$E(y^k) = \frac{\partial^k}{\partial t^k} M_y(t) \Big|_{t=0}$   $f(y)$

$= e^{\frac{1}{2}k^2} + \int_0^\infty y^k \sin(2\pi \log(y)) \, dy$

$= e^{\frac{1}{2}k^2} + \int_{-\infty}^\infty e^{ku} \sin(2\pi u) \, du$     let $u = \log y$

$= e^{\frac{1}{2}k^2} + 0$   $f(u)$     $du = \frac{1}{y} dy$

$= e^{\frac{1}{2}k^2}$

$-2$

This may again from the symmetric property of
$\sin(\cdot)$, or by using the geometric expansion of
$\sin(\cdot)$ such that shows the integral is 0.
However, due to the limit time and I couldn't
recall the formula for $\sin(\cdot)$ expansion, I cannot
perform this.

$\sin 2\pi u = \sin 2\pi (u-k)$

$\forall k \in \mathbb{N}$.

this is from

$\frac{\partial^k}{\partial t^k} \int_0^\infty e^{ty} \sin(2\pi \log(y)) f(y) \, dy \Big|_{t=0}$

$= \int_0^\infty e^{ty} y^k \sin(2\pi \log(y)) f(y) \, dy \Big|_{t=0}$

$= \int_0^\infty y^k \sin(2\pi \log(y)) f(y) \, dy$

$= E(y^k \sin(2\pi \log y))$

MGF
is not
finite in
this case!

or this question can be done by first
transform $y$ to $\log(y)$, then do the
computation, then the answer maybe
easier to get.

2. Let $X_1, X_2, \ldots, X_n$ be independent identically distributed (iid) random variables with common probability density function (pdf)

$$f_X(x|\theta) = \begin{cases} \frac{2x}{\theta} \exp\left\{-\frac{x^2}{\theta}\right\} & \text{if } x > 0 \\ 0 & \text{otherwise,} \end{cases}$$

where $\theta > 0$ is an unknown parameter.

(a) *(2 points)* Show that $T = \sum_{i=1}^n X_i^2$ is a complete sufficient statistic for $\theta$.

(b) *(4 points)* Determine the uniform minimum variance unbiased estimator (UMVUE) of $\theta^2$.

(c) *(4 points)* Determine the uniformly most powerful (UMP) rejection region for testing $H_0 : \theta \geq \theta_0$ versus $H_1 : \theta < \theta_0$, where $\theta_0$ is a known positive number. State your rejection region in terms of $T$ and a quantile from a named probability distribution.

Problem 2

(a) $f(x) = \prod_{i=1}^{n} \frac{2x_i}{\theta} \exp\left\{-\frac{x_i^2}{\theta}\right\} I(x_i > 0)$

$= \exp\left\{\sum_{i=1}^{n} \log(2x_i) - n\log(\theta) - \sum_{i=1}^{n} x_i^2/\theta\right\}$

$\in$ exponential family with $\eta(\theta) = \frac{1}{\theta}$, $T(x) = \sum_{i=1}^{n} x_i^2$

Since $\theta > 0$, $\frac{-1}{\theta}$ spans $\mathbb{R}$, $\exists$ open set, thus $T(x)$ is sufficient & complete for $\frac{-1}{\theta}$

Since $\frac{-1}{\theta}$ and $\theta$ is one-to-one, then $T(x)$ is sufficient & complete for $\theta$.

* can also use factorization theorem show $\sum x_i^2$ is minimal sufficient.

(b) $E(X^2) = \int_0^\infty \frac{2x^3}{\theta} \exp\left\{-\frac{x^2}{\theta}\right\} dx$   let $u = x^2$, $du = 2x\,dx$

$= \int_0^\infty \frac{u}{\theta} \exp\left\{-\frac{u}{\theta}\right\} du$

$= \theta$.

→ Recall $x \sim$ Gamma$(2, \frac{\theta}{\smile})$ ← Scale

$f_X(x) = \frac{1}{\Gamma(2)\theta^2} x^2 \exp\left(-\frac{x}{\theta}\right)$

→ first recognize, $Y = X^2 \sim \exp(\theta)$

Then $Var(Y) = \theta^2$

Thus, $\frac{1}{n-1}\sum_{i=1}^{n}(Y_i - \bar{Y})^2$ is an unbiased estimator of $\theta^2$.

$= \frac{1}{n-1}\sum_{i=1}^{n}(X_i^2 - \bar{X^2})^2$

Then by Rao-Blackwell & Lehmann Scheffé theorem.

UMVUE $= E[$unbiased estimator $|$ sufficient, complete statistics$]$

$= E\left[\frac{1}{n-1}\sum_{i=1}^{n}(X_i^2 - \bar{X^2})^2 \mid \sum_{i=1}^{n} X_i^2\right]$

$= E\left[\frac{1}{n-1}\sum_{i=1}^{n}(X_i^4 - 2X_i^2\bar{X^2} + (\bar{X^2})^2 \mid \sum_{i=1}^{n} X_i^2\right]$

please simplify

Let $Y = X^2$, $X = \sqrt{Y}$ (one-to-one since $X > 0$)   $J = \frac{\partial}{\partial y}\sqrt{y} = \frac{1}{2\sqrt{y}}$

$f_Y(y) = f_X(y) \cdot J$

$= \frac{2\sqrt{y}}{\theta} \cdot \exp\left\{-\frac{y}{\theta}\right\} \cdot \frac{1}{2\sqrt{y}}$

$= \frac{1}{\theta} \exp\left\{-\frac{y}{\theta}\right\}$

$\sim \exp(\theta)$   scale, mean parameter

(c) Recall from part (b) we recognize $Y = X^2 \sim \exp(\theta)$, then, a natural choice of statistics is $\frac{1}{n}\sum_{i=1}^{n} X_i^2$

Likelihood-ratio:  $\dfrac{L(\hat{\theta_0})}{L(\hat{\theta_0})} = \dfrac{(\frac{1}{\theta})^n \exp\left\{-\frac{\sum y_i}{\theta}\right\}}{(\frac{1}{\theta_0})^n \exp\left\{-\frac{\sum y_i}{\theta_0}\right\}}$   where $\hat{\theta_0}$ is the MLE estimator of $\theta$ under $H_0$

$\hat{\theta}$ is the MLE estimator of $\theta$ under $H_1$

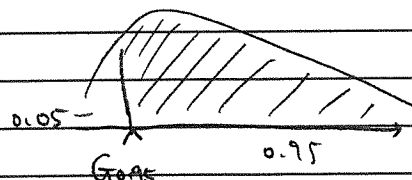why is monotonic in $\sum y_i = \sum x_i^2 = T(x)$, thus by Karlin Rubin's theorem, we reject when

Ratio is large, that is when $\frac{1}{n}\sum x_i^2 = \hat{\theta}$ is small (due to the $H_0$ direction)

Since $\frac{1}{n}\sum x_i^2 \sim$ Gamma$(n, \frac{\theta}{n})$, reject when

$\frac{1}{n}\sum x_i^2 < G_{(n,\frac{\theta}{n}),0.95}$

↑ 5th quantile for Gamma. $\alpha$ th

0.05

0.95

$G_{0.05}$

This is a UMP test.

3. The Pima Indians Diabetes data contains information of <u>768 women</u> from a population near Phoenix, Arizona, USA. The data contain medical records from Pima-Indian women at least 21 years of age. Below is reproduced the list of <u>5 attribute variables</u> and the class variable (response):

- 'glucose': plasma glucose concentration
- 'pregnancy': number of times pregnant
- 'bmi': body mass index (weight in kg/squared height in m$^2$)
- 'pedigree': diabetes pedigree function
- 'age': age in years
- 'diabetes': class variable (0 or 1)

The 'diabetes' variable is an indicator, with (1) indicating a positive test for diabetes between 1 and 5 years from the examination and (0) meaning a negative test.

We first fit a logistic regression model with 'diabetes' as the response variable and all the above attribute variables as predictors. The R code is produced below

```
model =  glm(diabetes~glucose+pregnancy+bmi+pedigree+age,
             data = PimaIndiansDiabetes2, family = binomial(link = "logit"))
summary(model)
```

and we have the following output:

```
Call:
glm(formula = diabetes ~ glucose + pregnancy + bmi + pedigree +
    age, family = binomial(link = "logit"), data = PimaIndiansDiabetes2)

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -8.673124   0.689980 -12.570  < 2e-16 ***
glucose      0.032893   0.003403   9.667  < 2e-16 ***
pregnancy    0.119458   0.031653   3.774 0.000161 ***
bmi          0.079550   0.013810   5.760 8.4e-09 ***
pedigree     0.891487   0.292239   3.051 0.002284 **
```

```
age              0.012230   0.009085   1.346 0.178247
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 993.48  on 767  degrees of freedom
Residual deviance: 732.51  on 762  degrees of freedom
AIC: 744.51

Number of Fisher Scoring iterations: 5
```

Then the code

anova(model)

gives the output here:

```
Analysis of Deviance Table

Model: binomial, link: logit

Response: diabetes

Terms added sequentially (first to last)
```

| | Df | Deviance | Resid. Df | Resid. Dev | Pr(>Chi) | |
|---|---|---|---|---|---|---|
| NULL | | | 767 | 993.48 | | |
| glucose | 1 | 184.764 | 766 | 808.72 | < 2.2e-16 | *** |
| pregnancy | 1 | 23.770 | 765 | 784.95 | 1.086e-06 | *** |
| bmi | 1 | 40.825 | 764 | 744.12 | 1.665e-10 | *** |
| pedigree | 1 | 9.819 | 763 | 734.31 | 0.001727 | ** |
| age | 1 | 1.797 | 762 | 732.51 | 0.180067 | |

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(a) *(3 points)* Write down the fitted model formula.

(b) *(4 points)* Conduct a hypothesis test at the 5% significance level that the coefficients for the two predictors 'pedigree' and 'age' are both zero. Give the numerical value of the test statistic and specify the rejection region.

(c) *(3 points)* We refit the model with the predictor 'age' being dropped and the summary output is below:

```
Call:
glm(formula = diabetes ~ glucose + pregnancy + bmi + pedigree,
    family = binomial(link = "logit"), data = PimaIndiansDiabetes2)

Coefficients:
             Estimate Std. Error z value Pr(>|z|)
(Intercept) -8.415851   0.656908 -12.811  < 2e-16 ***
glucose      0.033826   0.003345  10.112  < 2e-16 ***
pregnancy    0.141926   0.027105   5.236 1.64e-07 ***
bmi          0.078097   0.013771   5.671 1.42e-08 ***
pedigree     0.901294   0.291696   3.090    0.002 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 993.48  on 767  degrees of freedom
Residual deviance: 734.31  on 763  degrees of freedom
AIC: 744.31

Number of Fisher Scoring iterations: 5
```

Based on the fitted model, suppose that a woman in the study has a probability of 0.4 of having a positive diabetes test between 1 and 5 years from the examination. If her body mass index, i.e., the predictor 'bmi', is actually one unit less while the values of the other predictors remain the same, give the probability of this woman having a positive diabetes test.

Problem 3

−0.5 for not giving values of betas

(b) First write out the model:

$$\ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 \text{glucose} + \beta_2 \text{bmi} + \beta_3 \text{pedigree} + \beta_4 \text{age} + \underbrace{+ \beta_5 \text{pregnancy}}{} + \varepsilon$$

↓ intercept

↑ not correct −0.5

$H_0: \beta_3 = \beta_4 = 0$ $\qquad H_1: \text{not } H_0$

Recall the LRT test under GLM:

$$T = \frac{D_{M_0} - D_M}{\phi} \overset{H_0}{\sim} \chi^2_{p-q} \qquad \text{when sample size } n = 768 \text{ large}$$

$D_M = 732.51 \qquad \phi = 1 \qquad q = 762$

$D_{M_0} = 732.51 + 9.819 + 1.797 \qquad p = 762 + 2$

$$T = \frac{9.819 + 1.797}{1}$$

4

$$= 11.616 \overset{H_0}{\sim} \chi^2_2$$

$$> \chi^2_{2, 0.05} = 5.991$$

reject $H_0$ when $T > \chi^2_{2, 0.05} = 5.991$

In our case, we reject it, at least one of $\beta_3, \beta_4$ matters.

(a). Let $P_i = \text{Prob}(\text{Diabetes} \mid \text{other variables})$:

$$\log\left(\frac{P_i}{1-P_i}\right) = \beta_0 + \beta_1 \text{glucose}_i + \beta_2 \text{Preg}_i + \beta_3 \text{bmi}_i + \beta_4 \text{pedi}_i + \beta_5 \text{age}_i + \varepsilon_i$$

$i = 1, \ldots, 768$

(c). Origin $p = 0.4$

$$\log\left(\frac{P}{1-P}\right) = \log\left(\frac{0.4}{0.6}\right) = -0.405$$

with 1 unit decrease in bmi:

$$\Rightarrow \log\left(\frac{P}{1-P}\right)_{new} = -0.405 - 0.078097$$

$$= -0.48356$$

$$\Rightarrow$$

$$P_{new} = \frac{1}{1 + e^{0.48356}}$$

$$= 0.381$$

thus the true prob is 0.381.

3

4. Consider a study in which we count the number of defects in 1-square-inch specimens of fabrics. We have two different types of fabrics, A and B. Define the response variable $y$ as the count of defects and a binary variable $x$ that takes the values 0 or 1 when type is A or B respectively. The data is shown below.

| Specimen ID ($i$) | Type ($x_i$) | Count ($y_i$) |
|---|---|---|
| 1 | 0 | 1 |
| 2 | 0 | 0 |
| 3 | 0 | 0 |
| 4 | 0 | 1 |
| 5 | 1 | 2 |
| 6 | 1 | 1 |
| 7 | 1 | 2 |
| 8 | 1 | 3 |

Consider the following two models:

Model I:   $y \sim Poisson$, two types have two different means, $\lambda_A$ and $\lambda_B$

Model II:   $y \sim Poisson$, the mean is modeled as $\lambda = \beta_0 + \beta_1 x$.

Answer the following questions.

(a) *(2 points)* Write the mathematical forms of the two models described above.

(b) *(3 points)* Interpret the parameters in each model and find the relationship between $(\lambda_A, \lambda_B)$ and $(\beta_0, \beta_1)$.

(c) *(3 points)* Write the log-likelihood for model I. Using the data, find the Maximum Likelihood Estimator (MLE) of $(\lambda_A, \lambda_B)$. Give numerical value of the MLE.

(d) *(2 points)* Using the data, find the MLE of $(\beta_0, \beta_1)$ and give its numerical value.

Problem 4

(a) $\lambda = \beta_0 + \beta_1 X_i + \varepsilon_i$  $i = 1, \ldots, 8$. *    $Y \sim Poisson(\lambda)$

model 2    $f_y = \dfrac{\lambda^y e^{-\lambda}}{y!}$

model 1: mixed model depends on $I(x)$.

$\lambda = \lambda_A (1-X_i) + \lambda_B X_i$

$\quad = \lambda_A (1-X_i) + \lambda_B (X_i) + \varepsilon_i$

$\quad = \lambda_A + X_i (\lambda_B - \lambda_A) + \varepsilon_i$

· * This should be something Identical to the poisson model in ST704, which it in GLM with link $\ln(\cdot)$. But in that case, $\ln(\lambda)$ is linear, since we already have $\lambda$ linear, maybe do not need to log again?

(b). In model 2, $\beta_0$ is the expected mean when $X_i = 0$ (Type A)

$\quad\quad \beta_1$ is the expected change in mean when $X_i$ change from type A to type B.

In model 1, I'm fitting a linear model here using $\lambda_A$ and $\lambda_B$, and by comparing with model 2, we see $\beta_0 \equiv \lambda_A$, $\beta_1 \equiv (\lambda_B - \lambda_A)$

* However again, I'm not sure about the model, thus interpretation & relationships varies.

$\lambda_A$ & $\lambda_B$?

please see next page for part (a), (b) discussion.

c). $L(y) = \prod\limits_{i=1}^{4} \dfrac{\lambda_A^{y_i} e^{-\lambda_A}}{y_i!} \prod\limits_{i=5}^{8} \dfrac{\lambda_B^{y_i} e^{-\lambda_B}}{y_i!}$

$\ell(y) = C + \sum\limits_{i=1}^{4} y_i \log(\lambda_A) - 4\lambda_A + \sum\limits_{i=5}^{8} y_i \log(\lambda_B) - 4\lambda_B$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad C$ - constant invariant of $\lambda_A$, $\lambda_B$.

$\Rightarrow \dfrac{\partial \ell}{\partial \lambda_A} = \dfrac{\sum y_i}{\lambda_A} - 4 \overset{set}{=} 0$    $\dfrac{\partial^2 \ell}{\partial \lambda_A^2} = -\dfrac{\sum y_i}{\lambda_A^2} < 0$  concave.

$\Rightarrow \quad \hat{\lambda}_A = \dfrac{\sum y_i}{4} = \bar{y}$  $i = 1,2,3,4$    $\dfrac{\partial^2 \ell}{\partial \lambda_B^2} = -\dfrac{\sum y_i}{\lambda_B^2} < 0$

$\Rightarrow \dfrac{\partial \ell}{\partial \lambda_B} = \dfrac{\sum y_i}{\lambda_B} - 4 = 0$    Thus maximizers.

$\quad\quad\quad \hat{\lambda}_B = \bar{y}$  $i = 5,6,7,8$

$(\hat{\lambda}_A, \hat{\lambda}_B) = \left( \dfrac{2}{4}, \dfrac{2+1+2+3}{4} \right)$

$\quad\quad\quad\quad = \left( \dfrac{1}{2}, 2 \right)$

(d) by definition $\lambda_i = \beta_0 + \beta_1 X$.

$\ell(y) = \sum\limits_{i=1}^{8} y_i \log(\beta_0 + \beta_1 X_i) - 8(\beta_0 + \beta_1 X) + C$.

$\dfrac{\partial \ell}{\partial \beta_0} = \sum\limits^{8} \dfrac{y_i}{\beta_0 + \beta_1 X_i} - 8 \overset{set}{=} 0$

$\quad\quad\quad\quad\quad \beta_0 + \beta_1 \bar{X} = \bar{y}$

$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{X} = \bar{y} = \dfrac{5}{4}$    $\bar{X} = 0.5$  $\bar{y} = \dfrac{5}{4}$

$\dfrac{\partial \ell}{\partial \beta_1} = \sum\limits^{8} \dfrac{y_i X_i}{\beta_0 + \beta_1 X_i} - \sum X_i \overset{set}{=} 0$    $\sum (X_i - \bar{X})^2 = 0.5^2 \cdot 8 = 2$

$\quad\quad\quad \sum \beta_0 + \beta_1 X_i = \dfrac{\sum y_i X_i}{\sum X_i}$    $\sum(X_i - \bar{X})(y_i - \bar{y}) = (0.5 \cdot (1 - \dfrac{5}{4}) + 0.5 (-\dfrac{5}{4}))2$

result as SLR → $\hat{\beta}_1 = \dfrac{\sum (X_i - \bar{X})(y_i - \bar{y})}{\sum (X_i - \bar{X})^2}$    due to the time limit,    $+ (0.5(2 - \dfrac{5}{4}))2 + 0.5(1 - \dfrac{5}{4}) + 0.5(3 - \dfrac{5}{4})$

$\quad\quad = \dfrac{0}{2}$    I cannot recalculate these   $= (-\dfrac{1}{8} - \dfrac{5}{8})2 + (\dfrac{3}{8} \cdot 2) - \dfrac{1}{8} + \dfrac{7}{8}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ numbers, but the expressions

$\quad\quad = 0$    should be sound.    $= \dfrac{-6}{4} + \dfrac{3}{4} + (-\dfrac{1}{8}) + \dfrac{7}{8}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad = \dfrac{-6}{4} + \dfrac{3}{4} + \dfrac{3}{4}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\Rightarrow = 0 ?$

Problem4 continuous

(a)-(b). Given poisson distribution

$Y_i \sim$ Poisson $(\lambda) \Rightarrow f_y(y) = \frac{\lambda^y e^{-\lambda}}{y!} = P(y=y)$, $\ln(P(y=y)) = y\log(\lambda) - \lambda - \log(y!)$

And under GLM, we normally write $g(E(y|X)) = \beta_0 + \beta_1 X + \varepsilon$ where $g(\cdot)$ link function,

and the natural link function for poisson is $\log(\cdot)$ link, thus.

we normally have $\ln(\lambda) = \beta_0 + \beta_1 X + \varepsilon$

In this sense $\ln(\lambda_i) = \ln(\lambda_A)(1-X_i) + \ln(\lambda_B) X_i + \varepsilon_i$

$$= \ln(\lambda_A) + (\ln(\frac{\lambda_B}{\lambda_A}) X_i + \varepsilon_i \qquad model 1$$

In this case, $\ln(\lambda_A)$ is the expected log-mean when $X=0 \Rightarrow$ Type A.

$\ln(\frac{\lambda_B}{\lambda_A})$ is the expected log-mean change when $X \uparrow 1 \Rightarrow$ Type A to Type B

Similarly to model 2, but since already have $\lambda = \beta_0 + \beta_1 X$,

we may use identity link, that is

$$\lambda = \beta_0 + \beta_1 X + \varepsilon$$

Interpretation then consists with part (b).

5. An education researcher receives a grant to test the effectiveness of a new textbook (treat=1) compared to the standard textbook (treat=0). She randomly chooses $J$ schools in a large district to participate. For each school, exactly one class is chosen at random to participate in the trial, so we use a single index $j$ to refer to the school/class. All students within a class use the same textbook. Finally, the treatment is assigned at random.

Let $Y_{ij}$ denote a testing outcome (difference in test scores, after treatment minus before treatment) for student $i$ in classroom $j$ and treat$_j \in \{0, 1\}$ denotes treatment that is applied at the classroom level. The researcher assumes the following model:

$$Y_{ij} = \beta_0 + \beta_1 \cdot \text{treat}_j + v_{0j} + v_{1j} \cdot \text{treat}_j + \epsilon_{ij}$$

where:

$$v_{0j} \sim \mathcal{N}(0, \tau_0^2) \quad \text{(random intercept for classroom)}$$
$$v_{1j} \sim \mathcal{N}(0, \tau_1^2) \quad \text{(random offset for treat=1)}$$
$$\text{Cov}(v_{0j}, v_{1j}) = \tau_{01} \quad \text{(covariance between intercept and offset)}$$
$$\epsilon_{ij} \sim \text{i.i.d. } \mathcal{N}(0, \sigma^2) \quad \text{(residual error)}$$

The $\beta$ values are fixed effects, and we consider $\beta_1$ of primary interest.

(a) *(2 points)* From the model, what is var$(Y_{ij})$ for a randomly selected student in a random classroom that has been assigned the new textbook (treat=1)?

Below is partial R lme4 output from an analysis of the researcher's data, with 100 schools/classrooms (50 using the old text, 50 new), and classroom sizes ranging from 10 to 30 (average class size ~20).

```
> model <- lmer(y ~ treat + (1 + treat | class_id), data = df)
> summary(model)
Linear mixed model fit by REML ['lmerMod']
Formula: y ~ treat + (1 + treat | class_id)
   Data: df
```

```
Random effects:
 Groups    Name          Variance Std.Dev. Corr
 class_id (Intercept) 21.60    4.647
          treat         17.18    4.145    0.01
 Residual               15.85    3.982
Number of obs: 2066, groups:  class_id, 100

Fixed effects:
             Estimate Std. Error t value
(Intercept)  49.7791    0.6838   72.798
treat         2.4224    1.1132    2.176
```

(b) *(2 points)* For the null hypothesis that the treatment has no effect versus a two-sided alternative, would you reject the null at the 0.05 significance level? Provide your reasoning. (You may assume that the effective degrees of freedom is large enough that a standard normal approximation may be used in place of a $t$ distribution).

(c) *(2 points)* Suppose the researcher took the same data and simply averaged scores within each class, and used these average scores to perform an unequal variance two sample $t$-test of the 50 classes with the old text vs. the 50 with the new text. Would this be a valid approach, in the sense of controlling type I error? Explain.

(d) *(4 points)* Another researcher wants to replicate these results in a different city, where the average classroom size is 100, but they can afford to conduct the trial in only 50 classrooms (25 with the new text, 25 old). Using the output above, which study (first or second) do you think will have greater statistical power? Explain.

[HINT: Following part (c), think of observing a single class mean $\bar{Y}_j$ for each class, with $J/2$ classes per treatment group, and assume all classes are the same size $n$. Let $\bar{\bar{Y}}_{\text{treat}=k}$ denote the grand mean within treatment group $k$. The power of the t-test will be mainly determined by $\text{var}(\bar{\bar{Y}}_{\text{treat}=1} - \bar{\bar{Y}}_{\text{treat}=0})$, with small values providing greater power. From the output, you should be able to formulate $\text{var}(\bar{\bar{Y}}_{\text{treat}=1} - \bar{\bar{Y}}_{\text{treat}=0})$ in terms of $J$, $n$, and plugging in other estimates above, and you may assume $\tau_{01} = 0$. The first study has $J = 100$, $n = 20$, and the second study has $J = 50$, $n = 100$. Which study has smaller $\text{var}(\bar{\bar{Y}}_{\text{treat}=1} - \bar{\bar{Y}}_{\text{treat}=0})$?]

(a). $Var(Y_{ij})$

$= Var(\beta_0 + \beta_1 \cdot treat_j + V_{0j} + V_{ij} \cdot trt_j + \varepsilon_{ij})$

$= Var(V_{0j} + V_{ij} + \varepsilon_{ij})$   at trt $= 1$, drop fixed effects.

$= \sigma^2 + Var(V_{0j} + V_{ij})$   since $\varepsilon_{ij}$ should be indep from $V_{0j}, V_{ij}$

$= \sigma^2 + Var(V_{0j}) + Var(V_{ij}) + 2Cov(V_{0j}, V_{ij})$

if want numerical results $\rightarrow$   $= \sigma^2 + \zeta_0^2 + \zeta_1^2 + 2\zeta_{01}$ ✓

$= 15.85 + 21.6 + 17.18 + 0.02 = 54.65$

(b). $H_0: \beta_1 = 0$     $H_a: \beta_1 \neq 0$

$Z = \frac{\hat{\beta_1}}{SE(\hat{\beta_1})} = \frac{2.4224}{1.1132} = 2.176 > Z_{0.975} = 1.96$ ✓

Thus, we reject $H_0$ at $0.05$. ✓

(c). The approach itself is valid, that we are doing a 2 sample $t$-test.

however, under the unequal variance test, the df need to be adjusted using satterthwaite

✓ approximation. If the df is not been adjusted by the approx. method, the Type I error will be unpredicted.    good

To comparing the Type I error of doing such test with our current test, we recall the type I error is =

$$Prob(reject\ H_0 \mid H_0\ is\ correct) = \alpha$$

Thus when SE is larger, we tend to fail to reject $H_0$.

When comparing the variance of mean, the overall variance tend to be smaller than fitting a full model   (as $var(\bar{x}) = \frac{1}{n} var(x)$)

Thus are more likely to reject $H_0$, thus $\alpha$ may increase, not so optimal.

Also, testing on average loses information (indeed).   Actually, not much loss

(d). Now we have less $J$ but more $n$ in 2nd study.

for $J=100$, $n=20$,             for $J=50$, $n=100$

③ $Var(\bar{Y}_{..trt} - \bar{Y}_{..control})$       similarly,

$= Var(\bar{V}_{0.} + \bar{V}_{1.} + \bar{\varepsilon}_{..} - \bar{V}_{0.} - \bar{\varepsilon}_{..})$      $Var(\bar{Y}_{..trt} - \bar{Y}_{..control})$

$= Var(\bar{V}_{1.})$                           $= \frac{1}{50}\zeta_1^2$

$= \frac{1}{100}\zeta_1^2$                      $= 0.3436$

$\approx 0.1718$

Thus, first study have smaller variance,     More complicated tradeoff thus first test has more power.     between $J$ and $n$