

Semester 1 2017,
Example Exam 2

Instructions

The exam consists of 4 questions, 1-4. Each question has four items, a-d.

Within each question:

Item (a) carries a weight of 8 marks.

Item (b) carries a weight of 7 marks.

Item (c) carries a weight of 6 marks.

Item (d) carries a weight of 4 marks.

The total marks in the exam are 100.

Answer ALL questions in the spaces provided.

If more space is required, use the back of the PREVIOUS page.

Show all your working and include sketches where appropriate.

Work written in the formulae and tables booklet will NOT be marked.

Question 1:

The weight of a component (in kilograms) is continuous random variable, X with support $[0, 1]$ and pdf

$$f(x) = \begin{cases} Ke^{-x}, & x \in [0, 1], \\ 0, & \text{otherwise,} \end{cases}$$

where K is some constant.

(a) Find K .

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

$$= \int_0^1 ke^{-x} dx = 1$$

$$= [-ke^{-x}]_0^1 = 1$$

$$= -ke^{-1} + ke^0 = 1$$

$$k(1 - e^{-1}) = 1$$

$$k = (1 - e^{-1})^{-1}$$

$$\approx 1.582$$

(b) Calculate an expression for the cdf of X and plot it.

$$F(x) = \int f(x) dx$$

$$= -ke^{-x} + C$$

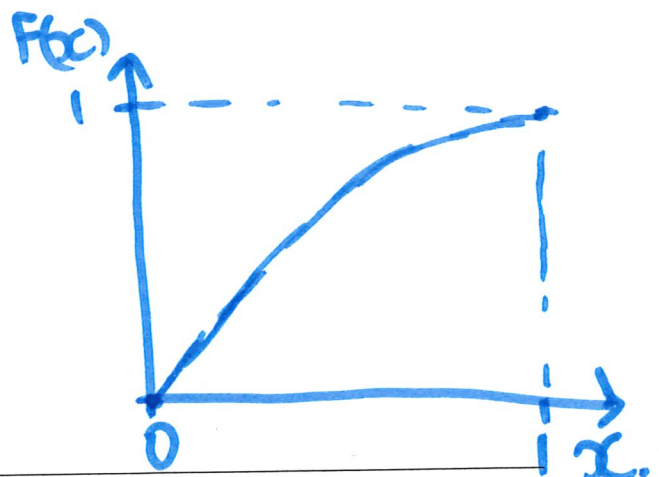
$$F(0) = 0$$

$$-ke^0 + C = 0$$

$$\therefore C = k$$

$$F(x) = -1.582e^{-x} + 1.582$$

$$= 1.582(1 - e^{-x})$$



(c) Assume now that 7 such independent random variables are measured and let N denote the number of those components having weight greater than $1/2$. What is $P(N \geq 2)$?

$$\begin{aligned} \Pr(X > 0.5) &= 1 - \Pr(X \leq 0.5) \\ &= 1 - F(1/2) \\ &= 1 - 0.622 \\ &= 0.378 \end{aligned}$$

$$\begin{aligned} N &\sim \text{bin}(7, 0.378) \\ \Pr(N \geq 2) &= 1 - \Pr(N < 2) = 1 - [\Pr(N=0) + \Pr(N=1)] \\ &= 1 - [0.622^7 + 7 \times 0.622^6 \times 0.378] \\ &= 1 - 0.189 = 0.811 \end{aligned}$$

(d) Continuing on the previous item, for each of the 7 independent random variables, you are paid \$5 if it has a weight less than $1/2$ and \$3 otherwise. What is the variance of your total earnings?

$$N \sim \text{bin}(7, 0.378)$$

$$\text{var}(N) = np(1-p)$$

key

$$\begin{aligned} *E &= 3N + 5(7-N) \\ &= 3N + 35 - 5N \\ &= -2N + 35 \end{aligned}$$

$$\begin{aligned} n &= 7 \\ p &= 0.378 \end{aligned}$$

$$\text{var}(E) = \text{var}(-2N + 35)$$

$$= \text{var}(-2N)$$

$$= (-2)^2 \text{var}(N)$$

$$= 4 \times 7 \times 0.378 \times 0.622$$

$$= 6.58$$

**this is very challenging*

Question 2:

A homework assignment consists of 10 questions. The chance of succeeding in a question is $p = 0.8$ and the success/failure of each question is independent of the others.

There are two marking schemes:

Scheme A: All questions are marked and each is worth 10%.

Scheme B: Only a randomly selected subset of size 5 is marked and each is worth 20%.

In both schemes a question is either correct (success) or wrong (failure). Let X denote the percent grade under scheme A and let Y denote the percent grade under scheme B.

(a) What is the mean of X ? How about the mean of Y ? In this respect which marking scheme is better? Or are the schemes equivalent?

$$\begin{aligned}
 X &= 10(Q_1 + Q_2 + Q_3 + \dots + Q_{10}) & E(X) &= E(Y) \\
 E(X) &= 10E(Q_1 + Q_2 + \dots + Q_{10}) & & \therefore \text{Schemes} \\
 &= 10(0.8 + 0.8 + \dots + 0.8) & & \text{are} \\
 &= 80 & & \text{equivalent}
 \end{aligned}$$

$$\begin{aligned}
 Y &= 20(Q_1^* + Q_2^* + \dots + Q_5^*) \\
 E(Y) &= 20(0.8 + 0.8 + 0.8 + 0.8 + 0.8) \\
 &= 80
 \end{aligned}$$

(b) What is $P(X = 100)$? How about $P(Y = 100)$? Now in this respect which marking scheme is better for the student?

$$N_x \sim \text{bin}(10, 0.8) \quad N_y \sim \text{bin}(5, 0.8)$$

$$\Pr(N_x = 10) = 0.8^{10} = 0.107$$

$$\Pr(N_y = 5) = 0.8^5 = 0.328$$

scheme Y is better in this regard.

(c) What is the variance of X ? How about the variance of Y ? In this respect, which marking scheme is better? In what way is it better?

$$\text{var}(N) = n p (1-p)$$

$$\begin{aligned} \text{var}(N_x) &= 10 \times 0.8 \times 0.2 \\ &= 1.6 \end{aligned}$$

$$\begin{aligned} \text{var}(N_y) &= 5 \times 0.8 \times 0.2 \\ &= 0.8 \end{aligned}$$

$$X = 0.1 N_x$$

$$Y = 0.2 N_y$$

$$\begin{aligned} \text{var}(X) &= 0.1^2 \times 1.6 \\ &= 0.016 \end{aligned}$$

$$\text{var}(Y) = 0.2^2 \times 0.8$$

$$= 0.032$$

can't tell which is better. more likely to get high & low marks under scheme Y.

(d) Failure occurs if the percent grade is less than 50%. Compare the chance of failure in both schemes. Which scheme is better for the student?

$$\Pr(N_x < 5) \quad \text{or} \quad \Pr(N_y < 3)$$

$$\begin{aligned} \Pr(N_y \leq 2) &= \Pr(N_y=0) + \Pr(N_y=1) + \Pr(N_y=2) \\ &= 0.2^5 + \binom{5}{1} 0.2^4 \times 0.8 + \binom{5}{2} 0.2^3 \times 0.8^2 \end{aligned}$$

$$= 0.058$$

$$\begin{aligned} \Pr(N_x < 5) &= \Pr(N_x=0) + \Pr(N_x=1) + \Pr(N_x=2) + \Pr(N_x=3) \\ &\quad + \Pr(N_x=4) \end{aligned}$$

⋮

$$= 0.0064$$

∴ scheme X is better in terms of failures.

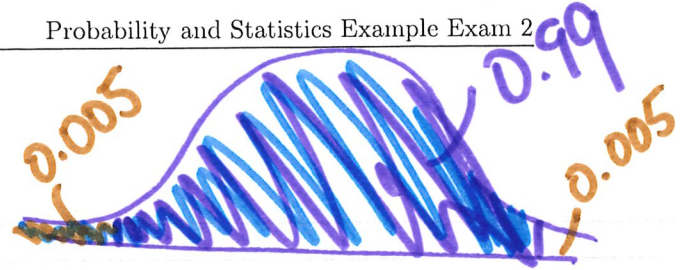
Question 3:

You observe the following output:

```
In [5]: UnequalVarianceTTest(data1,data2)
Out[5]: Two sample t-test (unequal variance)
-----
Population details:
parameter of interest: Mean difference
value under h_0: 0
point estimate: -1.4516996384754872
95% confidence interval: (-5.91045104330129,3.0070517663503153)

Test summary:
outcome with 95% confidence: fail to reject h_0
two-sided p-value: 0.47361548935859044 (not significant)

Details:
number of observations: [7,8]
t-statistic: -0.7524144182475475
degrees of freedom: 7.902492459938804
empirical standard error: 1.9293883839395962
```



$\bar{x} - \bar{y}$

$$\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

df

(a) Present a 99% confidence interval for the difference in means under the assumption of unequal variances.

$$\bar{x} - \bar{y} \pm t \times \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

$$t_{8;0.995} = 3.355$$

$$-1.452 \pm 3.355 \times 1.929$$

$$-1.45 \pm 6.48$$

$$(-7.93, 5.03)$$

(b) You now know that the sample standard deviation of "data2" is 2.028. What is the sample standard deviation of "data1"?

$$\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} = 1.93^2$$

$$= \frac{S_1^2}{7} + \frac{2.028^2}{8} = 1.93^2$$

$$S_1^2 = 7 \left(1.93^2 - \frac{2.028^2}{8} \right) = 22.48$$

$$S_1 = 4.74$$

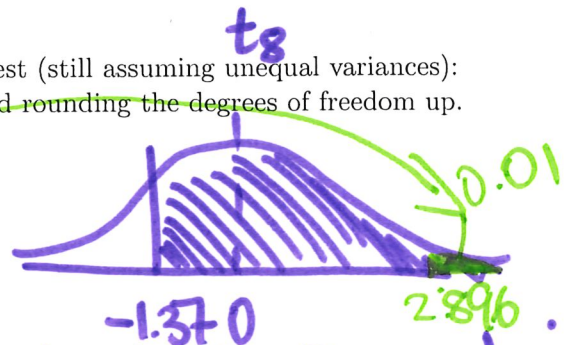
(c) Use the same data to carry out a one-sided hypothesis test (still assuming unequal variances):
 $H_0: \mu_1 = \mu_2 + 1.2$ vs. $H_1: \mu_1 > \mu_2 + 1.2$ with $\alpha = 0.01$ and rounding the degrees of freedom up.

$$H_0: \mu_1 - \mu_2 = 1.2 \quad H_1: \mu_1 - \mu_2 > 1.2$$

$$t_s = \frac{\bar{x} - \bar{y} - 1.2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$= \frac{-1.45 - 1.2}{1.93}$$

$$= -1.37$$



p-value > 0.5 ∴ retain H_0
 because $-1.37 < 2.896$,
 retain H_0
 there's no evidence
 that $\mu_1 - \mu_2 > 1.2$

(d) Carry out the same test using the assumption of equal variances.

$$t_s = \frac{\bar{x} - \bar{y} - 1.2}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

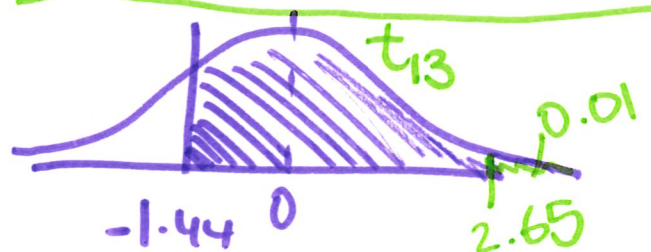
$$= \frac{-1.45 - 1.2}{3.55 \sqrt{\frac{1}{7} + \frac{1}{8}}}$$

$$= \frac{-2.65}{1.84} = -1.44$$

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

$$= \sqrt{\frac{6 \times 4.74^2 + 7 \times 2.028^2}{13}}$$

$$= \sqrt{12.58} = 3.55$$



because $-1.44 < 2.65$
 retain H_0

Question 4:

Assume the hypothetical situation where you have observations:

$$(x_1, y_1) = (10, 10),$$

$$(x_2, y_2) = (20, 22),$$

$$(x_3, y_3) = (30, 32),$$

$$(x_4, y_4) = (40, 40).$$

(a) Calculate the sample mean and sample variance for both (x_1, x_2, x_3, x_4) and (y_1, y_2, y_3, y_4) .

$$\bar{x} = (x_1 + x_2 + x_3 + x_4) / 4 \quad | \quad \bar{y} = (10 + 22 + 32 + 40) / 4$$

$$= (10 + 20 + 30 + 40) / 4 \quad | \quad = 26$$

$$= 25 \quad | \quad \text{var}(y) = (10^2 + 22^2 + 32^2 + 40^2) / 4$$

$$\text{var}(x) = (\sum x_i^2 - n\bar{x}^2) / (n-1) \quad | \quad = (3208 - 2704) / 3 \quad \frac{-426}{3}$$

$$= (10^2 + 20^2 + 30^2 + 40^2 - 4 \times 25^2) / 3 \quad | \quad = 504$$

$$= (3000 - 2500) / 3 \quad | \quad \frac{500}{3}$$

$$= 500 / 3$$

(b) Compute (by hand) the least squares estimates, for the model,

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, 2, 3, 4,$$

where $\epsilon_i \sim N(0, \sigma^2)$.

$$\sum y_i x_i = 10 \times 10 + 20 \times 22 + 30 \times 32 + 40 \times 40$$

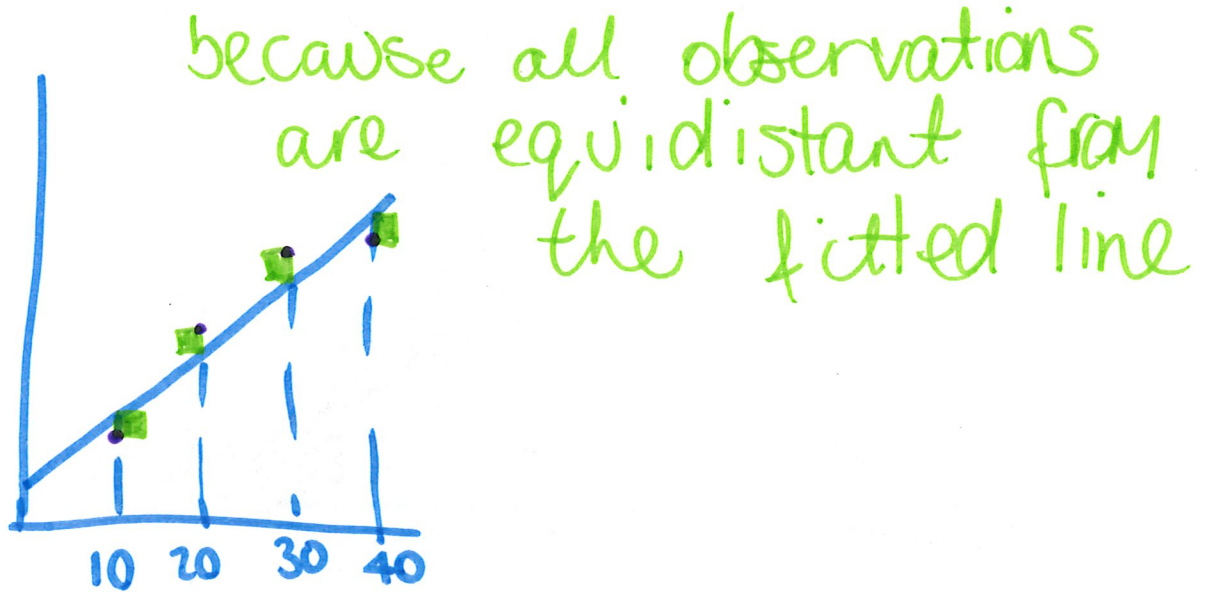
$$= 3100$$

$$\hat{\beta}_1 = \frac{\sum y_i x_i - (\sum y_i)(\sum x_i) / n}{\sum x_i^2 - (\sum x_i)^2 / n} = \frac{3100 - \frac{104 \times 100}{4}}{3000 - \frac{100^2}{4}}$$

$$= \frac{500}{500} = 1$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 26 - 1 \times 25 = 1$$

(c) Sketch a plot of the regression line on a plane containing the ⁴ three data points. Explain why this line minimises the sum of squares agrees with your answer to (b).



(d) Assume now that after estimating the parameters, you use the model,

$$Y = \hat{\beta}_0 + \hat{\beta}_1 x + \epsilon,$$

$$\epsilon \sim \text{unif}(\quad)$$

where $\hat{\beta}_0$ and $\hat{\beta}_1$ are according to your answer in (b) but ϵ is uniformly distributed (as opposed to Normally distributed) with mean 0 and variance 1. Sketch the cdf of Y under this model when $x = 20$.

$$\mu = 0 = \frac{a+b}{2}$$

$$\sigma^2 = (b-a)^2 = 1$$

$$-a = b$$

$$(b-a)^2 = \frac{3}{2}$$

$$a = -\sqrt{\frac{3}{2}}$$

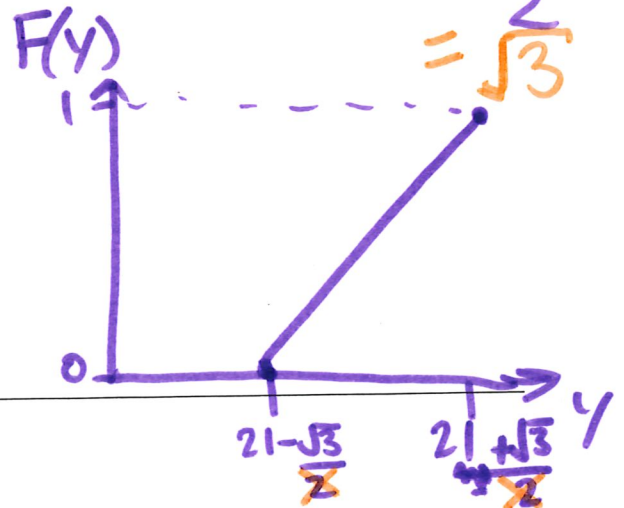
$$b-a = \sqrt{\frac{3}{2}} \Rightarrow 2b = \sqrt{\frac{3}{2}}$$

$$b = \sqrt{\frac{3}{2}}$$

$$Y = \hat{\beta}_0 + \hat{\beta}_1 x + \epsilon$$

$$= 1 + 1 \times 20 + \epsilon$$

$$= 21 + \epsilon$$



$$Y \sim \text{unif}\left(21 - \frac{\sqrt{3}}{2}, 21 + \frac{\sqrt{3}}{2}\right)$$