

2061551
P1

1. (i) In study one, if a person's unhealthy diet increase by 1 unit, the odds of developing CHD increases by 21.1%. ~~the odds of developing CHD increases by 21.1%.~~
~~Expressed multiple times, the 95% and 95% confidence interval of~~
the odds of developing CHD ranging from decreasing 17.5% to increasing 77.8%. Since the 95% CI ~~includes~~ includes 1, this means the effect is not statically significant. //

66

2/2

(ii)

$$OR_2 = \frac{\text{odds of failure on treatment}}{\text{odds of failure on control}} = \frac{\frac{45}{198}}{\frac{35}{221}} = \frac{221}{154} \approx 1.4351$$

$$V(\log(OR_2)) = \frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d} = \frac{1}{45} + \frac{1}{35} + \frac{1}{198} + \frac{1}{221} \approx 0.0604 //$$

✓ 2/2

(iii)

$$W_i = \frac{b_i \times c_i}{n_i}, \text{ where } b_i \text{ is number of people with risk factor present and unhealthy,}$$

c_i is number of healthy people with risk factor absent.

n_i is number of people

where b_i is number of healthy people with risk factor present in group i ,

c_i is number of diseased people with risk factor absent in group i ,

n_i is number of all people being investigated in group i

2/2

(iv)

$$W_2 = \frac{198 \times 35}{499} = \frac{6930}{499} \approx 13.888 //$$

✓ 1/1

2561551

p 2

$$(V.) \quad OR_{MH} = \frac{\sum_{i=1}^3 w_i \times OR_i}{\sum_{i=1}^3 w_i}$$

$$OR_2 = 1.435$$

$$V(\log OR_2) = 0.0604$$

$$w_2 = 13.888$$

$$\Rightarrow OR_{MH} = \frac{23.7 \times 1.21 + 13.888 \times 1.435 + 33.193 \times 1.276}{23.7 + 13.888 + 33.193}$$

$$= 1.285$$

Study	OR	$V(\log OR)$	w
1	1.21	0.0384	23.7
2	1.435	0.0604	13.888
3	1.276	0.0255	33.193

$$(Vi) \quad V[\log(OR_{MH})] = \frac{\sum_{i=1}^3 w_i^2 V_i}{(\sum_{i=1}^3 w_i)^2}, \text{ where } w_i = \frac{b_i \times c_i}{n_i} \text{ (same as part ii)} \\ v_i = \frac{1}{a_i} + \frac{1}{b_i} + \frac{1}{c_i} + \frac{1}{d_i}$$

a_i : ~~number of~~ diseased people with risk factor present in group i
 b_i : number of healthy people with risk factor present in group i
 c_i : number of diseased people with risk factor ~~present~~ absent in group i
 d_i : number of healthy people with risk factor absent in group i
 n_i : number of people being investigated in group i.

(vii)

$$95\% \text{ CI for } OR_{MH}: \exp[\log(OR_{MH}) \pm 1.96 \sqrt{V(\log(OR_{MH}))}]$$

$$V[\log(OR_{MH})] = \frac{23.7^2 \times 0.0384 + 13.888^2 \times 0.0604 + 33.193^2 \times 0.0255}{(23.7 + 13.888 + 33.193)^2} = 0.0122$$

$$\Rightarrow 95\% \text{ CI for } OR_{MH}: (\exp(\log(1.285) \pm 1.96 \times \sqrt{0.0122}))$$

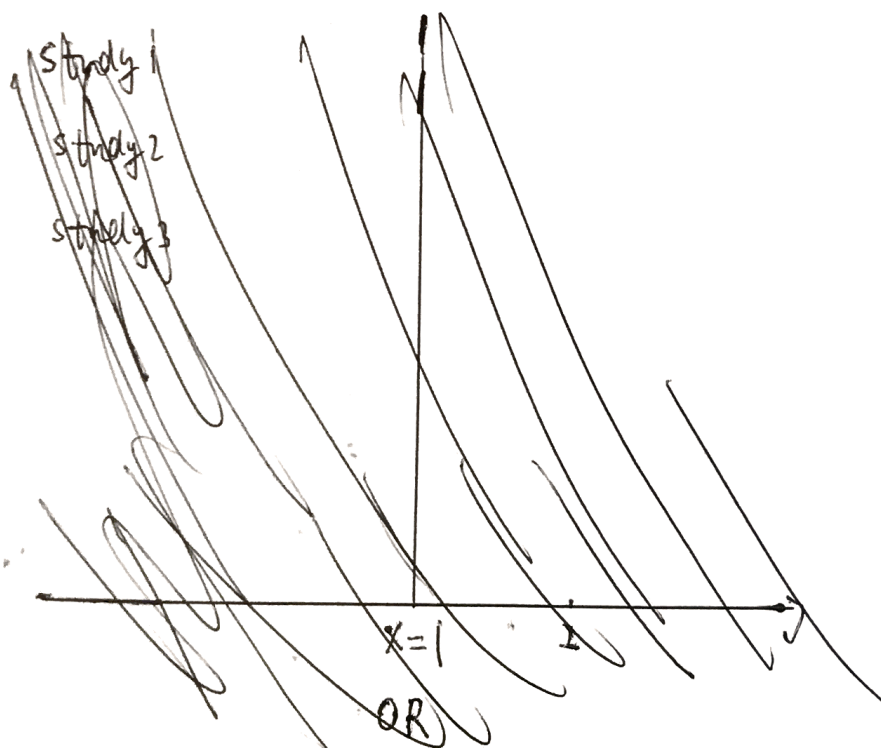
$$1.115$$

$$(1.029, 1.193)$$

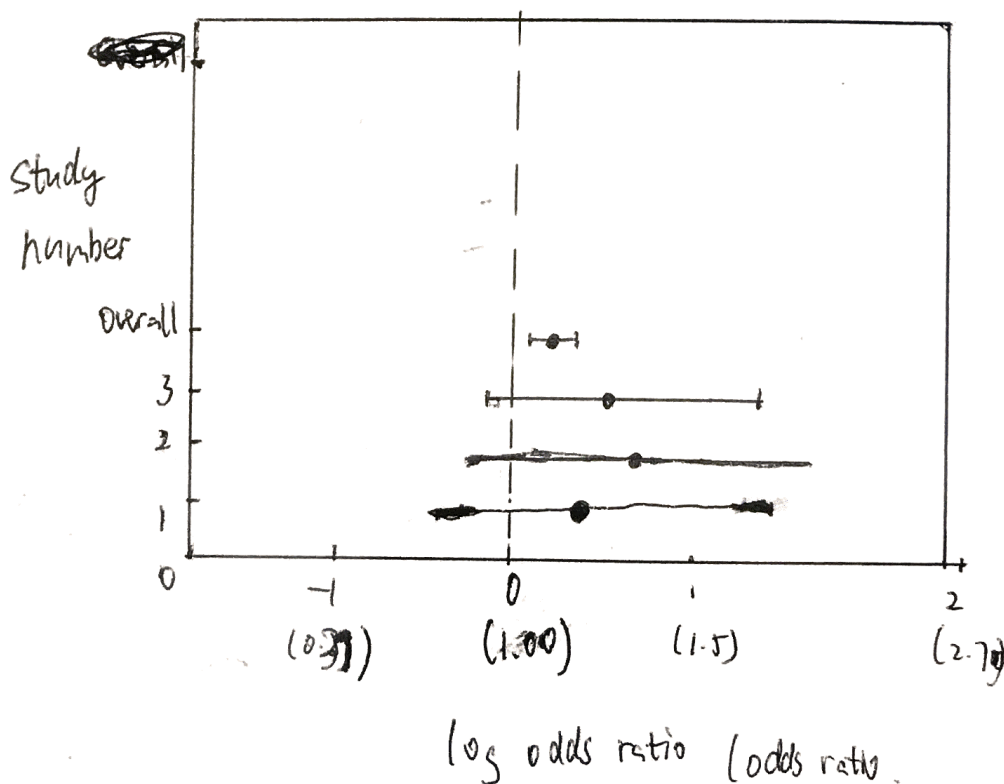
3.5/4

(viii) If a person's unhealthy diet increase by 1 unit, the odds of developing CHD increases by ~~11.5%~~, with ~~95%~~ percent confidence interval of increase ranging from 2.9% to 18.3%. ~~Since 95% of Mantel-Haenszel Z test~~ ~~does not contain 1~~, the effect is statically significant. 1/2

(ix)



Graphic display of results of meta analysis



✓

5/5

2561551

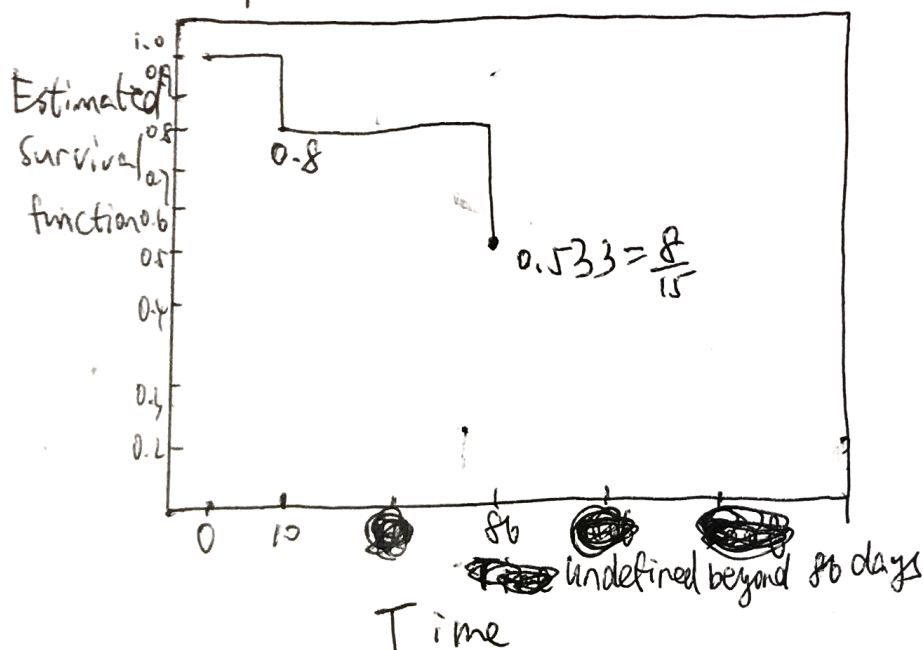
P4

2.(i)

Event time: $t(i)$	At risk at time $t(i)$: $r(i)$	surviving beyond $t(i)$: $s(i)$	Estimated survival function at $t(i)$: $\hat{S}(t(i))$
0	5	5	$\frac{5}{5} = 1$
10	5	4	$1 \times \frac{4}{5} = 0.8$
86	3	2	$0.8 \times \frac{2}{3} = 0.533$ $= \frac{8}{15}$

✓ 3/3

(ii) Kaplan-Meier Estimator illustration



2/2

(iii) $P(\text{survive at least 30 days}) = 0.8$ //

1/1

(iv) $S(t) = \exp(-\int_0^t h(u) du) = \exp(-H(t)) \Leftrightarrow H(t) = -\ln(S(t))$

at time point $(t=0)$: $H(0) = -\ln(1) = 0$

at time point $(t=10)$ $H(10) = -\ln(0.8) \approx 0.223$

$t = 86$, $H(86) = -\ln(\frac{8}{15}) \approx 0.629$ //

✓ 2/2

256(551)

P5

$$(V.) \hat{H}_{NA}(t) = \sum_{j=1}^i \frac{d_j}{r_j}, \text{ where } d_j = r_j - s_j$$

$$\Rightarrow \text{at } t=0, \hat{H}_{NA}(0) = \frac{5-5}{5} = 0$$

$$\text{at } t=10, \hat{H}_{NA}(10) = \frac{5-4}{5} + 0 = \frac{1}{5}$$

$$\text{at } t=86, \hat{H}_{NA}(86) = 0 + \frac{1}{5} + \frac{3-2}{3} = \frac{8}{15}$$

✓ 2/2

(vi) parametric analysis means fitting specific families of distribution to the survival data, for example, log-normal, gamma etc.

Nelson and Aalen cumulative hazard estimator is a non-parametric estimator.

✓ 2/2

(b)

(i) H_0 : The survival functions for Group 1 and Group 2 are the same

H_1 : The survival functions for Group 1 and Group 2 are different

2/2

(ii)	d_{1i}	d_{2i}	r_{1i}	r_{2i}	$E_{1i} = \frac{r_{1i}}{r_i} d_i$	$E_{2i} = \frac{r_{2i}}{r_i} d_i$
42	1	0	4	4	$\frac{4}{8} \times 1 = \frac{1}{2}$	$\frac{4}{8} \times 1 = \frac{1}{2}$
227	1	0	2	4	$\frac{2}{6} \times 1 = \frac{1}{3}$	$\frac{4}{6} \times 1 = \frac{2}{3}$
1996	0	1	0	2	0	$\frac{2}{2} \times 1 = 1$
2878	0	1	0	1	0	$\frac{2}{2} \times 1 = 1$
			$O_1 = 2$	$O_2 = 2$	$E_1 = \frac{5}{6}$	$E_2 = \frac{19}{6}$

$$\chi^2_{LR} = \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} = 2.063 < \chi^2_1 = 3.84$$

5/5

2561551

P. 8

(iii) Since $\chi^2_{LR} = 2.063 < \chi^2 = 3.84$

\Rightarrow we fail to reject H_0 and conclude that the survival function for Group 1 and Group 2 are the same. //

2/2

(c)

(i) hazard rate for age: $\exp(\hat{\beta}_1) = 1.03$

Interpretation: Increase the age by 1 unit, the hazard ~~rate~~ rate increase by 3%

2/3

~~(ii) $h_0(t)$ is an arbitrary baseline hazard function, evaluated when all the covariates are 0.
 z_i is the covariate vector =~~

$h_0(t)$ is an arbitrary baseline hazard function, when age and t_5 are all at default value 0.

z_i is the covariate vector = $\begin{pmatrix} \text{age (in years)} \\ t_5 \text{ (continuous)} \end{pmatrix}$

β is the set of parameters: $\beta^T = (\beta_1: \text{age}, \beta_2: t_5)$

pf let $z_1 = \begin{pmatrix} 12 \\ 1.26 \end{pmatrix}$ $z_2 = \begin{pmatrix} 13 \\ 1.48 \end{pmatrix}$

2.5/3

$$\frac{h(t, z_1)}{h(t, z_2)} = \frac{h_0(t) \exp(z_1^T \beta)}{h_0(t) \exp(z_2^T \beta)} = \exp((z_1 - z_2)^T \beta) \text{ irrelevant of time}$$

256/551

p7

3. (a)(i)

COVID-19

	present	absent	Total
White	385	331,464	331,849
ethnic minority	64	16,685	16,749
Total	449	348,149	348,598

COVID-19

	present	absent	Total
ethnic minority	64 ^A	16,685 ^B	16,749
White	385 ^C	331,464 ^D	331,849
Total	449	348,149	348,598

(ii)

$$P_1 = P(\text{Disease} | \text{White}) \approx 0.00116$$

$$P_2 = P(\text{Disease} | \text{minority}) \approx 0.00382$$

$$AR = 2.648 \text{ per } 1000 \text{ ppl}$$

Interpretation: the absolute change in risk is an increase of 2.648 people per 1000 people if you have the risk factor ethnic minority.

$$(iii) \quad RR = \frac{\frac{64}{64+16685}}{\frac{385}{331464+385}} = 3.294$$

$$\text{Var}(\ln RR) = \frac{1}{A} - \frac{1}{A+B} + \frac{1}{C} - \frac{1}{C+D} = 0.018$$

95% CI in log scale: (0.929, 1.455)

95% of RR: (2.532, 4.284)

There is almost 3-fold increased risk if you are ethnic minorities, and the effect is statically significant as the 95% CI interval does not include 1.

575

2561551

p8

(iv)

$$OR = \frac{\frac{64}{16685}}{\frac{385}{331464}} \approx 3.302$$

⇒ there is almost 3-folds increased risk of getting COVID-19 if you are ethnic minorities compared to white subjects per 1000 ppl. In 1000 people, the number of people getting COVID-19 of ethnic minority is almost 3 times higher than white subjects. 2/2

(v) No, they are not.
~~OR measures the ratio difference, while~~

$$RR = \frac{\frac{A}{A+B}}{\frac{C}{C+D}}, \quad OR = \frac{\frac{A}{B}}{\frac{C}{D}}$$

odds ratio measure how risk effect affect the outcome, whereas relative risk measure the whole group, meaning that OR is not stable as RR. 0/2

(vi) ~~Poisson relative risk, logistic regression~~

Study A measures relative risks, Study B measures odds ratios.

Because in Poisson regression relative risk is natural in Poisson model, while odds ratio fits well is logistic regression. 2/2

2561551

P9

$$(b) (i) \text{ sensitivity} = \frac{\text{number of } \overset{\text{diseased people screen positive}}{\cancel{\text{not screen positive}}}}{\text{Total number of diseased ppl}}$$

If a screen test has low sensitivity, which means given the total number of diseased ppl, the test cannot identify the disease well, it only identify a few people who have disease, there are many people with disease, and the screen test cannot identify. ✓✓

$$(ii) \text{ sensitivity} = \frac{214}{372} = 0.575$$

$$\text{specificity} = \frac{5291}{5312} = 0.996$$

$$(FNR) \text{ false negative rate} = 1 - 0.575 = 0.425$$

$$(FPR) \text{ false positive rate} = 0.004$$

A sensitivity of 0.575 means 57.5% of diseased people are correctly identified as having COVID-19.

A false negative rate of 0.425 means 42.5% of diseased people are not correctly identified as having COVID-19.

A specificity of 0.996 means 99.6% of healthy people are correctly identified as ^{not} having COVID-19.

A false positive rate of 0.004 means 0.4% of healthy people are not correctly identified as having COVID-19.

5/5

501551

P10

I recommend strategy **A**,
Consider two tests in row. in plan B
COVID-19

	Y	N
positive flow test	214 { 123 Correct 91 wrong }	21 { 20.9 Correct 0.1 wrong }
negative flow	158 { 91 Correct 67 wrong }	5281 { 5273 Correct 21 wrong }

Sensitivity (B) = $\frac{123}{372} = 0.331 = P(\text{two screen positive} | \text{diseased})$ **Ans B**

Specificity (B) = $\frac{5273}{5312} = 0.99 = P(\text{two screen negative} | \text{absent})$

For plan A.

since two event independent, sensitivity (A) = 0.575

specificity (A) = 0.996

→ choose plan A.

X 0/9