

# QE sheet 2

(note, Q worded a bit like it wants you to find the cdf...)

4a) We're interested in the distribution of  $N$ , where  $N$  is the number of trials such that  $\sum_{i=1}^N X_i = 1$ , with each  $X_i$  iid  $\sim$  Bernoulli( $p$ ).

$N=1$  w.p.  $p$  ie success on trial 1  
 $N=2$  w.p.  $(1-p)p$  ie.  $X_1 \bar{X}_2$  and so on.  
 So  $P(N=n) = (1-p)^{n-1} \cdot p$  for all  $n \in \mathbb{N}^+$

~~it actually~~  
~~I think I~~  
~~misread Q,~~  
~~we actually~~  
~~want~~  
 $(1-p)^n$

b)  $P(N=n) = \begin{cases} 0 & \text{if } n=1 \\ 1-p^n - (1-p)^n & \text{? otherwise} \end{cases}$   
 for  $n \in \mathbb{N}^+$

I don't follow...

For  $n \geq 2$ , in order to not meet the condition, we require either all  $\bar{X}_i$  or all  $X_i$

7a) i)  $f_X : P(X=0) = 0.3$   
 $P(X=1) = 0.7$

$f_Y : P(Y=0) = 0.4$   
 $P(Y=1) = 0.1$

$P(Y=2) = 0.5$

Marginal pmf is pmf regardless of others' values

ii. ~~pmf of Y~~

		y		
		0	1	2
x	0	$1/3$	0	$2/3$
	1	$3/7$	$1/7$	$3/7$

numbers in cells are  $P(Y=y | X=x)$

b) Not independent.  $X \perp\!\!\!\perp Y \Leftrightarrow$  the distribution of  $Y$  does not depend on the value of  $X$ , i.e.  $f_{Y|X=x}$  is the same for all  $x$ .

Clearly this isn't true, as shown above.

cor  $P(Y|X=x) =$

c)  $E[Y] = 0 \times 0.4 + 1 \times 0.1 + 2 \times 0.5 = 1.1$

$E[Y|X=1] = 1/7 + 2 \times 3/7 = 1$

$E[Y|X=0] = 4/3$

$P(Y)$ , in different notation! ie

and  $0.3 \times 4/3 + 0.7 \times 1 = 1.1$  as expected. thah.

d)  $E[Y|X] \stackrel{\text{def}}{=} \sum_{y \in \{0,1,2,3\}} y \cdot P(Y=y|X)$

Yes [This object is a function of the random variable  $X$ , and thus itself a random variable.]

9a)  $\text{Cov}(Z, aX + bY + c) \stackrel{\text{def}}{=} E[Z \cdot (aX + bY + c)] - E[Z]E[aX + bY + c]$

$= E[Z \cdot aX] - E[Z]E[aX] + E[Z \cdot bY] - E[Z]E[bY] + E[Z \cdot c] - E[Z]E[c]$  by linearity of  $E$

$= a \cdot (E[Z \cdot X] - E[Z]E[X]) + b \cdot (\dots) + c \cdot (E[Z] - E[Z])$   
by properties of  $E$  for constants, scaling, also from linearity

$\stackrel{\text{def}}{=} a \text{Cov}(Z, X) + b \text{Cov}(Z, Y)$  as required

b) As noted, let  $Z := aX + bY + c$ .

$\text{Var}(Z) \stackrel{\text{def}}{=} \text{Cov}(Z, Z) \equiv E[Z^2] - E[Z]^2$

$= a \text{Cov}(Z, X) + b \text{Cov}(Z, Y)$  from a)

$= a \cdot (a \text{Cov}(X, X) + b \text{Cov}(X, Y)) + b \cdot (b \text{Cov}(Y, Y) + a \text{Cov}(X, Y))$

by a) with appropriate variable subs. (I think this would work too)

$\stackrel{\text{def}}{=} a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$  as required

for  $*$ ,  $\text{Cov}(aX + bY + c, X) = E[X \cdot (aX + bY + c)] - E[X]E[aX + bY + c]$   
 $= E[X] \cdot E[bY + c] - E[X]E[bY + c] + \text{by linearity}$   
 $= aE[X^2] - aE[X]^2 = a \text{Var}(X)$   
and the same for the second term.

13. a)  $X \sim \text{Bern}(p)$

$$\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$$

a)  $E[\bar{X}_n] \stackrel{\text{def}}{=} E\left[\frac{1}{n} \sum_{i=1}^n X_i\right]$

$$= \frac{1}{n} E\left[\sum_{i=1}^n X_i\right]$$

$$= \frac{1}{n} \sum_{i=1}^n E[X_i]$$

$$= \frac{1}{n} \cdot \sum_{i=1}^n p$$

$$= p$$

} by linearity

as each  $X_i$  is iid taking 1 w.p.  $p$  and 0 otherwise

~~$\text{Var}(\bar{X}_n) \stackrel{\text{def}}{=} E[\bar{X}_n^2] - p^2$  from just above~~

$\text{Var}(\bar{X}_n) \stackrel{\text{def}}{=} \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right)$

and as each  $X_i$  is iid,  $\text{Var}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \cdot n \cdot \sigma^2$  where  
 $\sigma^2 := \text{Var}(X_i) = \text{Var}(X) \stackrel{\text{def}}{=} E[X^2] - E[X]^2 = p - p^2$   
 so  $\text{Var}(\bar{X}_n) = \frac{1}{n} \cdot p \cdot (1-p)$

b) The CLT states that  $\lim_{n \rightarrow \infty} Z_n \sim N(0, 1)$

where  $Z_n := \frac{\bar{X}_n - E(\bar{X}_n)}{\sqrt{\text{Var}(\bar{X}_n)}}$  is the standardised sample mean of size  $n$ .

i.e. it is normally distributed with mean 0 and variance 1.

Let  $L$  be the number of failed light bulbs.

$$P(15 \leq L \leq 22) = P(L \leq 22) - P(L \leq 14)$$

$$P(L \leq k) = P\left(\bar{X}_n \leq \frac{k}{100}\right) = P\left(Z \leq \frac{\frac{k}{100} - 0.2}{\sqrt{0.2 \cdot 0.8}}\right) \text{ where } Z \sim N(0, 1)$$

$$0.94 = 0.935$$



Let  $\hat{p}$  be the observed proportion of failed patients.

c)  $P(0.15 \leq \hat{p} \leq 0.22) = P(0.15 \leq \bar{X}_{100} \leq 0.22)$   
 where  $\mu = 0.2$ ,  $\sigma^2 = \frac{1}{100} \times 0.2 \times 0.8 = 0.0016$ ,  $\sigma = 0.04$   

$$= P\left(\frac{0.15 - 0.2}{0.04} \leq Z \leq \frac{0.22 - 0.2}{0.04}\right)$$
  

$$= P(Z \leq 0.5) - P(Z \leq -1.25) = 0.691 - 0.106$$
  

$$= 0.59$$

15 a)  $H_0: \bar{p}_x - \bar{p}_y = 0$   $H_1: \bar{p}_x - \bar{p}_y > 0$

I'd avoid the bars — usually used for sample mean! —  
 where  $\bar{p}_{x,y}$  is the <sup>population</sup> proportion of patients recovering who <sub>do, do not</sub> receive the drug.  
 In our sample, the  $z$  statistic

$$Z = \frac{\hat{p}_x - \hat{p}_y}{\sqrt{\text{Var}(\hat{p}_x - \hat{p}_y)}}$$

where  $\hat{p}_x = 0.75$   
 $\hat{p}_y = 0.65$

and  $\text{Var}(\hat{p}_x - \hat{p}_y)$  under  $H_0 = \frac{\hat{\pi}(1-\hat{\pi})}{100}$

$$\frac{\hat{\pi}(1-\hat{\pi})}{100} \cdot 2 \quad \text{with } \hat{\pi} = 0.7$$

By CLT and with a large iid sample,  $Z \sim N(0,1)$  given  $H_0$ .

For a 1-tailed test, the critical values at the 1% and 10% significance levels are 1.282 and 2.326 (satisfying  $P(Z \geq C_\alpha) = \alpha$ .)

As observed,  $Z = \frac{0.1}{\sqrt{0.0042}} = 1.543$ .

$1.543 < 2.326$ , so insufficient evidence to reject  $H_0$  at 1%  
 but  $1.543 > 1.282$ , so we can reject  $H_0$  at 10%.

b) The p-value tells us the smallest significance level  $\alpha$  at which we can reject  $H_0$ . In this case the p-value, i.e.  $P(Z > 1.543) = 6.14\%$ . So we could reject  $H_0$  at 10% but not 1%.

when you conclude falsely there's an effect, w.p.  
c) A Type I error is  $P(\text{reject } H_0 \mid H_0 \text{ is true}) = \alpha$ , the significance level. Type II is  $P(\text{fail to reject} \mid \text{false}) = \beta$ ,  
when you conclude falsely there's no effect, w.p.

where  $1 - \beta$  is the power of a test.

We'd probably prefer to avoid a Type I error, since prescribing an ineffective drug could be very harmful, but it depends a lot on side-effects, disease severity, alternatives, etc.

d) With  $n = 300$ ,  $z = \frac{0.1}{\sqrt{\frac{0.0014}{300}}} = 2.67$  and now the result is significant at 1%. This isn't especially surprising, more trial participants increases power, as the variance in the different recovery rates is smaller with our larger sample so we can be more confident the observed difference wasn't noise.

Great