

Question 1

76.2

1.
$$X = \begin{cases} 1, & p = \frac{1}{3} \\ 3, & p = \frac{1}{2} \\ 4, & p = \frac{1}{6} \end{cases}$$

a)
$$E(X) = \sum x_j \cdot f(x_j)$$

$$= 1 \cdot \frac{1}{3} + 3 \cdot \frac{1}{2} + 4 \cdot \frac{1}{6} = \frac{11}{6}$$

b)
$$\text{Var}(X) = \sum (x_i - \mu)^2 \cdot p_{X_i}$$

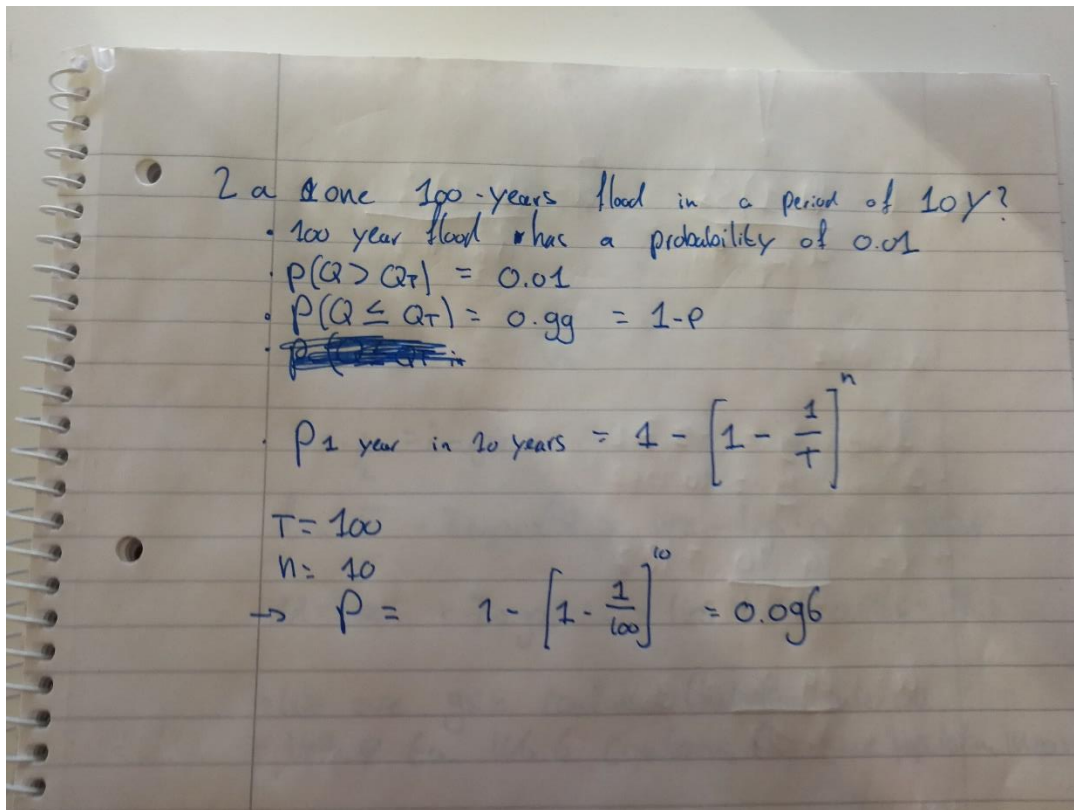
$$= \left(1 - \frac{11}{6}\right)^2 \cdot \frac{1}{3} + \left(3 - \frac{11}{6}\right)^2 \cdot \frac{1}{2} + \left(4 - \frac{11}{6}\right)^2 \cdot \frac{1}{6}$$

$$= \frac{149}{36} \approx 4.14$$

c. Mode = 3. Since $X=3$ has the highest probability

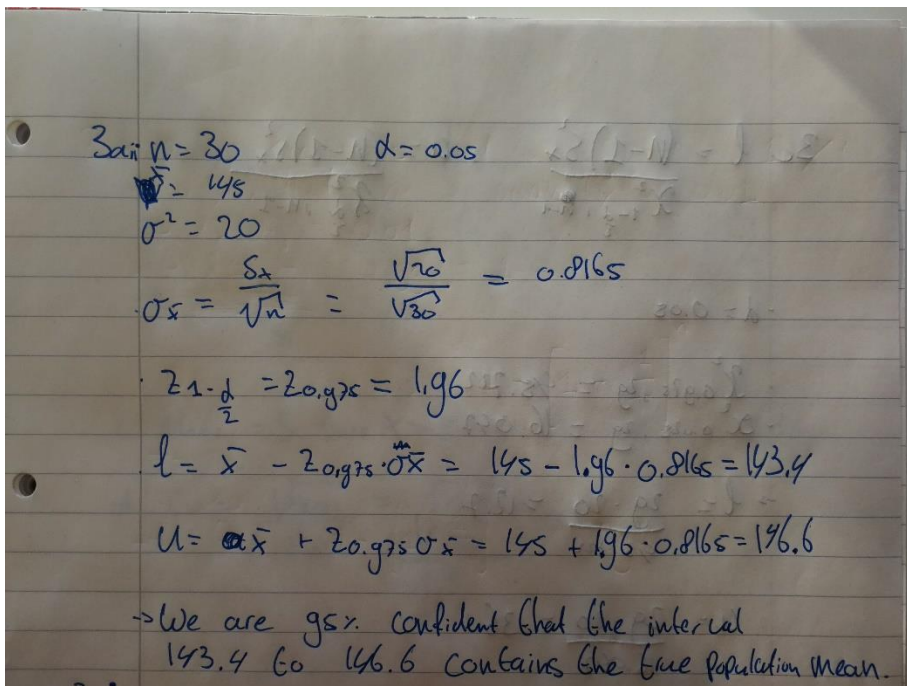
d)
$$Cv = \frac{\sqrt{\text{Var}(X)}}{E(X)} = \frac{\sqrt{\frac{149}{36}}}{\left(\frac{11}{6}\right)} \approx 1.11$$

Question 2



2b) The assumption of normality is violated in this analysis. The data have more extreme values than would be expected from a normal distribution. The use of a lognormal distribution (or another heavy-tail distribution) instead of a normal distribution can improve the analysis.

Question 3 (I did 3ai before 3ai, therefore the referring to 3aii in 3ai)



3ai

Same procedure as 3aii, but using t-value instead of z-value:

- $s_x = 0.8165$
- $t_{1-\frac{\alpha}{2}, n-1} = t_{0.975, 29} = 2.045$
- $l = \bar{x} - t_{0.975, 29} \cdot s_x = 143.3$
- $u = \bar{x} + t_{0.975, 29} \cdot s_x = 146.6$

→ we are 95% confident that the interval 143.3 to 146.6 contains the true population mean.

3b) The reason for the difference is that in part i we assume the true variance is not known and therefore assume a t-distribution with n-1 degrees of freedom instead of the z-value that we use in part ii with infinite degrees of freedom.

3c

$$l = \frac{(n-1)s_x^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2}, \quad u = \frac{(n-1)s_x^2}{\chi_{\frac{\alpha}{2}, n-1}^2}$$

- $\alpha = 0.05$
- $\chi_{0.975, 29}^2 = 45.722$
- $\chi_{0.025, 29}^2 = 16.047$

→ $l = \frac{29 \cdot 20}{45.722} = 12.7$

$u = \frac{29 \cdot 20}{16.047} = 36.1$

→ The 95% confidence intervals for the variance are found to be 12.7 to 36.1.

4a) In machine learning, a model is trained on the data in the training set. The best possible fit for this data is constructed. The prediction error is estimated with a loss function (i.e. the mean squared error). This error is the training error. Then the model is used to predict the outcome of the test set. The test set is NOT USED in order to construct the model. Again the prediction error (now for the test set) is estimated. This is the test error. The test error is a good estimation of how the model performs in practice (because the data is not used to construct a learner).

4b) We want to control model complexity because a too complex model might lead to overfitting. This means the model predicts the training data very well (a low training error), but performs poorly on the test data (high test error). And since the goal is to get a model with a low test error, we need to control the model complexity.

5a) We can fit a linear regression to the time series and then perform a t-test to test if the slope, beta, of the regression is significantly different from 0. The null hypothesis is that there is no trend. The alternative hypothesis is that there is a trend. The test statistic is:

$$t = \frac{\beta - 0}{S_{\beta}}, \text{ and the critical value is } t_{1-\alpha/2, n-2}$$

S_{β} is the standard deviation of beta. The null hypothesis is rejected if $|t| > t_{1-\alpha/2, n-2}$

5b) The times series is not a perfect sine or cosine wave, so it is not C. From the times series there seems to be a strong periodicity of around 5 s. $1/5 = 0.2$ hz. Graph A shows a peak at $f = 0.2$ hz, so graph A shows the Fourier transform of X_t .