

# DATA ANALYSIS PROJECT

RAVURU CHIDAKSH(200010046)

200010046@iitdh.ac.in

Department of Computer Science, IIT Dharwad

September 24, 2021

# Contents

<b>1</b>	<b>Approach towards Problem 2</b>	<b>3</b>
1.1	Calculation of mean and variance: . . . . .	3
1.2	Central Idea of the Approach: . . . . .	3
<b>2</b>	<b>Images explaining the method</b>	<b>4</b>
<b>3</b>	<b>Note:</b>	<b>8</b>

# 1 Approach towards Problem 2

We are given a dataset with 50000 samples in it, of the random variable  $Z$ . We also know that  $Z = X + 10Y$  where  $Y = \sum_{i=0}^k W_i$  and we know  $W_i$  are independent and identically distributed random variables. And we know  $k$  belongs to  $\{2,3,4\}$  and  $W_i$  belong to one of either Exponential distribution or Rayleigh distribution or Half-normal distribution.

## 1.1 Calculation of mean and variance:

From the given data set we can calculate mean and variance of the entire data, I took the help of pandas series objects and data frames to load the data from file and then took help of numpy inbuilt methods `np.mean` and `np.var` to calculate the mean and variance of the data.

## 1.2 Central Idea of the Approach:

Since we have multiple values of  $k$  and multiple distributions we need to find the best possible  $k$  and distribution for the given data which means that we need to find that  $k$  and distribution such that the original mean and variance of data satisfies the relation between mean and variance of one of the distributions so closely for some  $k$  in  $\{2,3,4\}$ .

So we calculate mean for some  $k$  in  $\{2,3,4\}$  and then find variance (say `var1`) assuming some distribution and then we know actual variance from the data (say `var2`) and we will check for  $\min(\text{var1} - \text{var2})$  for every distribution for all  $k$ .

## 2 Images explaining the method

Given ,  $Z = X + 10Y$

$X \rightarrow$  Uniform distribution

$Y = \sum_{i=1}^k W_i$  where  $k \in \{2, 3, 4\}$

Applying linearity of Expectation,

$$E[Z] = E[X] + 10 \times E[Y]$$
$$\Rightarrow E[Z] = E[X] + 10E[Y]$$

and  $E[X] = \frac{b+a}{2}$  (for uniform distribution with ends  $[a, b]$ )

$$= \frac{3+(-3)}{2}$$
$$= 0.$$
$$E[Y] = E\left[\sum_{i=1}^k W_i\right]$$

We know  $w_i$ 's are i.i.d

→ let  $w_i$ 's are random variables of  $W$  which follows one of rayleigh, exponential or half normal distributions,

$$\begin{aligned}\Rightarrow E[Y] &= E[w_1 + w_2 + \dots + w_k] \\ &= E[w_1] + E[w_2] + \dots + E[w_k] \\ &= k \cdot E[w] \quad (\text{all are } E[w])\end{aligned}$$

$$\therefore \boxed{E[Y] = k \cdot E[w]}$$

$$\Rightarrow E[Z] = 0 + 10 \cdot k \cdot E[w]$$

$$\Rightarrow \boxed{E[w] = \frac{E[Z]}{10k}} \rightarrow \textcircled{1}$$

Similarly for variance,

$$\text{var}(Z) = \text{var}(X + 10Y)$$

( $X, Y$  are independent because if  $X, Y$  are dependent it would make ~~some~~ all of  $w_i$ 's dependent which is a contradiction).

and hence  $X, Y$  are independent,

$$\text{var}(Z) = \text{var}(X) + 100 \cdot \text{var}(Y)$$

$$\downarrow$$
  

$$\frac{(b-a)^2}{12} \quad (b=3, a=-3)$$

$$\Rightarrow \text{var}(Z) = \frac{(3 - (-3))^2}{12} + 100 \cdot \text{var}(Y)$$

$$\Rightarrow \text{var}(z) = 3 + 100 \cdot \text{var}(y)$$

$$\Rightarrow \frac{\text{var}(Z) - 3}{100} = \text{var}(Y) \quad \text{--- } \textcircled{\times}$$

but  $\text{var}(y) = \text{var}(w_1 + w_2 + \dots + w_k)$  ( $\because w_i$ 's are iid)

$$= \text{var}(w_1) + \text{var}(w_2) + \dots + \text{var}(w_k)$$

$$= K g^2 \quad (\text{all are } \text{var}(w))$$

$$\therefore \frac{\text{var}(\bar{z}) - 3}{100 \cdot k} = \sigma^2(w) \quad \rightarrow (2)$$

Now as we know  $E[w]$ ,  $\sigma^2[w]$  from data



→ For some  $k$  in  $(2, 3, 4)$  say (2) we can assume some distribution (say exponential) and calculate  $\text{var}[w]$  using the relation between mean and variance for that particular distribution. Say that variance ~~be~~ be  $\sigma^2[w_{\text{new}}]$ .

by finding ratio between  $\sigma^2[w_{\text{new}}]$  to  $\sigma^2[w]$  for all 3 possible  $k$ 's and 3 possible distributions, the case in which the ratio is close to 1 the corresponding  $k$  and distribution is what we are looking for.

→ We can also take difference instead of ratio and check for which combination it's zero.

### 3 Note:

- I wrote the code in python in Google Colab , which i'm submitting with extension "ipynb". So, to run the entire file , please upload the file on drive and open it with google colabaratory.
- I worked with .xls file in the entire process which i'm going to zip it with this pdf file . Please upload that .xls file on Google colab before running the cells.