

Nikhil Sawlani

i) $\mu = 77.0$, $\sigma = 3.4$, Find probability.

a) 72.6, less than 72.6.

Soln $\rightarrow Z = \frac{X - \mu}{\sigma}$ X - value to be standardized
 μ - mean, σ - std deviation.

$$\rightarrow Z = \frac{72.6 - 77}{3.4} = -1.294, \text{ Z-value of } (-1.294) = \#$$

$$\rightarrow \text{pnorm}(-1.294) \quad \text{or by Z table}$$

$$\rightarrow + 0.09783 \quad \text{or} \quad 0.0985$$

b) greater than 88.5

$$\rightarrow 1 - Z(88.5)$$

$$\rightarrow 1 - Z\left(\frac{88.5 - 77}{3.4}\right) = \cancel{1 - Z(3.38)} \quad 1 - Z(3.38)$$

$$= 1 - 0.99964$$

$$\rightarrow \cancel{Z(-2.3)} = 0.00036$$

c) Between 81 & 84.

$$\therefore \frac{84 - 77}{3.4} = 2.05, \quad Z(2.05) \rightarrow 0.9798 \quad \text{--- (a)}$$

$$\& \quad \frac{81 - 77}{3.4} = 1.17, \quad Z(1.17) \rightarrow 0.8790 \quad \text{--- (b)}$$

$$\therefore a - b = 0.1008$$

d) Between 56 & 92.

$$\therefore \frac{92 - 77}{3.4} = \frac{92 - 77}{3.4} = 4.4, \quad Z(4.4) = \text{pnorm}(4.4) \quad \text{--- (a)}$$

$$\& \quad \frac{56 - 77}{3.4} = -6.17, \quad Z(-6.17) = \cancel{\text{pnorm}(-6.17)} \quad \text{--- (b)}$$

$$\therefore a - b = 0.9999946 \approx 1$$

2) $\mu = 145$, $sd = 12 = 6$

(a) → percentage between 130 & 160 data points.

$$z = \frac{X - \mu}{\sigma} = \frac{130 - 145}{12} = \frac{-15}{12} = -1.25$$

For negative value, $1 - z(-1.25) = 0.10565$ (a)

$$z = \frac{160 - 145}{12} = 1.25, \quad z(1.25) = 0.8944 \quad (b)$$

$$\therefore, b - a = 0.8944 - 0.1056 = 0.7888$$

(b) → less than 130 minutes.

$$\rightarrow \frac{130 - 145}{12} = \frac{-15}{12} = -1.25 = 0.10565$$

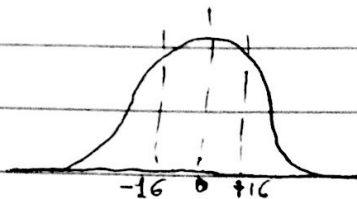
(3) a) For any question answer can be True or false.
For question 1, 2, ..., 15, probability of true or false is $\frac{1}{2}$ for True or $\frac{1}{2}$ for false.

for 5 questions → $\left(\frac{1}{2}\right)^5$

b) → for any 5 questions → $\left(\frac{1}{2}\right)^5$ 15cs

$$= \left(\frac{1}{2}\right)^5 \times \frac{15!}{10! \times 5!}$$

4) 68% between 35 & 42



$$Z = \frac{X - \mu}{\sigma}$$

Here, we have 68% area between 35 & 42 i.e. lowest point & highest point.

one standard deviation.

$$\therefore -1 = \frac{35 - \mu}{\sigma} \quad \& \quad 1 = \frac{42 - \mu}{\sigma}$$

$$-6 = 35 - \mu \quad \text{--- (a)} \quad \& \quad 6 = 42 - \mu \quad \text{--- (b)}$$

Now, multiply a & b.

$$-6 = 35 - \mu \quad \text{--- (a)}$$

$$6 = 42 - \mu \quad \text{--- (b)}$$

$$35 - \mu + 42 - \mu = 0$$

$$77 = 2\mu$$

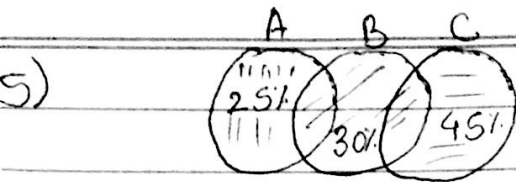
$$\mu = 38.5$$

Now, substituting μ value in Z formula --- (a)

$$-6 = 35 - 38.5$$

$$6 = 3.5$$

$$\therefore \text{Mean} = 38.5 \quad \& \quad \text{Sd} = 3.5$$



Here, we need to find what percent of A is set of B.

$$\therefore \frac{A}{B} = \frac{25}{30} = \frac{5}{6} = 0.833 \approx 83.3\%$$

6) $P(T/U) = 82\%$, $P(T \cap D) = 62\%$, $P(D \cap T) = ?$

$$P(D \cap T) = \frac{P(T \cap D)}{P(T/U)} = \frac{62\%}{82\%} = 75.60\%$$

7) All three heads →
 a) we have only one combination of such case.
 $HHH \rightarrow \therefore \frac{1}{8}$, as sample space is $2^3 = 8$ for 3 coins.

b) Exactly one head → $\frac{3}{8}$, (HTT, THT, TTH)

c) $P(H / (2 \text{ more heads})) = \frac{4}{8} = \frac{1}{2}$

8) → $\mu_1 = 1200$, $\mu_2 = 1800$
 $\text{Var}_1 = 90000$, $\text{Var}_2 = 160000$
 $X_1 = ?$, $X_2 = ?$, $\text{Sd}_1 = 300$, $\text{Sd}_2 = 400$
 $P(X_1) > P(X_2)$.

Here, to calculate probability of claims filed by 2nd agent is lower, so → $X_1 - X_2$

1) $N_1 = 1000$ \rightarrow given

$$\mu_1 = 270 \text{ seconds}$$

$$\sigma_1 = 100 \text{ seconds}$$

$$\text{Variance} = 10000$$

$$N_2 = ?$$

$$\mu_2 = ?$$

$$\sigma_2 = 10 \text{ seconds}$$

$$\text{Variance}_2 = 100$$

Soln \rightarrow

$$\text{Var} = npq$$

$$10000 = npq$$

Here, $np \sim \text{mean}$ so we can replace np by mean

$$10000 = 270q$$

$$\therefore q = 37.03$$

$$\& \quad p = \frac{\text{Mean or } np}{N_1}$$

$$P = \frac{270}{10000}$$

$$P = 0.027$$

for N_2 ,

$$V_2 = N_2 P_2 q_2$$

$$100 = N_2 (0.027)(37.03)$$

$$N_2 \approx 100$$

∴ difference between 2 random variable

$$\rightarrow \phi\left(\frac{\mu_1 - \mu_2}{\sqrt{\sigma_1^2 + \sigma_2^2}}\right)$$

we have given in question, Mean1, Mean2, sd1 & sd2 & Variance 1 & var2.

Ans

$$\therefore \phi\left(\frac{1200 - 1800}{\sqrt{90000 + 160000}}\right) = \frac{-600}{\sqrt{250000}} = \frac{-600}{500} = -1.2$$

Z value of -1.2 is \rightarrow $\text{pnorm}(-1.2) \rightarrow 0.1150697$

Answer \rightarrow 11.5 %

Q \rightarrow

10)

$$P(S) = 60\%$$

$$P(Sp) = 99\%$$

$$P(D) = 5\%$$

$$P(SP') = 1\% \text{ (No drug use)}$$

S-Sensitivity

Sp-Specificity

Let Employee test positive be E.

$$P(E/D) \text{ i.e. Sensitive} = 60\%$$

$$P(E/D') = \text{Specificity} = 99\% \text{ (all employee under test will be under 99\%)}$$

find, Employees test positive, that actually use drugs,

$$P(D/E) = \frac{P(E/D) P(D)}{P(E)}$$

$P(E) \rightarrow$ whole probability

\rightarrow Sensitive

$$= \frac{P(E/D) P(D)}{P(E/D) P(D) + P(E/D') P(D')}$$

\rightarrow Not Sensitive

$$= \frac{60 \times 5}{60 \times 5 + 1 \times 99} = 0.76$$

11) two bags solution

y = year from which bag came 1994.

E → yellow ~~color~~ color of bag

Find,

$$P(y/E) = \frac{P(E/y) P(E)}{P(E)}$$

$$= \frac{P(E/y) P(E)}{P(E/y) P(E) + P(E/y') P(E')} \quad \text{--- (a)}$$

So, we have possible outcomes from the years as yellow & green exist in both the years.

If we take yellow & green from 1994 then

$y \rightarrow 20\%$, $G \rightarrow 10\%$.

& for 1996, $y \rightarrow 14\%$, $G \rightarrow 20\%$.

putting value for 1994,

$$= \frac{20 \times 20 \times \frac{1}{2}}{20 \times 20 \times \frac{1}{2} + 14 \times 10 \times \frac{1}{2}} \quad \text{--- Here, } \frac{1}{2} \text{ is the probability for bag from any 2 bags. years.}$$

$$= \frac{400}{400 + 140} = \frac{400}{540}$$

= 0.74 is the probability that we have chances of yellow bag from 1994.

12.1 -

```
A<-read.csv("WMT_1.CSV")
```

```
A
```

```
mean(A$Close)
```

```
sd(A$Close)
```

```
View(A)
```

12.2 -

```
B<-read.csv("SHLD.CSV")
```

```
B
```

```
C<-data.frame(B)
```

```
class(C)
```

```
mean(B$Close)
```

```
sd(B$Close)
```

```
View(B)
```

12.3 -

The coefficient of variation(CV) talks about dispersion of data in any data sets.

CV function in R is a part of raster library.

Therefore,

```
library(raster)
```

```
cv(A$Close)
```

```
cv(b$Close)
```

12.4 -

CV helps to identify dispersion in data. Ex as in real life scenarios we can find CV values for any given share from historical data. Lesser the CV higher the reliability.

12.5 - A for Walmart, B for Kmart, C for DJIA

`cv(c$Close) = 4.6741236`

`cv(A$Close) = 2.938513`

`cv(B$Close) = 3.543339`

CV is minimum for A then B then C, which implies A has minimum dispersion and A has minimum risk, which makes A more attractive to any investor.

```
A<-read.csv("WMT_1.CSV")
```

A

```
mean(A$Close)*100
```

```
sd(A$Close)*100
```

13.1

```
install.packages('pastecs')
```

```
library('pastecs')
```

```
stat.desc(a$orders)
```

```
-----
```

```
library(e1071)
library(ggplot2)
a<-read.csv("Prob_Assignment_Dataset.CSV")
summary(a$orders)
skewness(a$orders)
kurtosis(a$orders)
plot(a$site,type="l")
plot(a$visits~a$orders)
boxplot(a$visits~(a$platform))
boxplot(a$orders~(a$platform))
boxplot(a$orders~a$new_customer)
boxplot(a$visits~a$new_customer)
```

Bivariate Analysis

14.1

```
ggplot(a, aes(x = factor(a$visits),a$orders)) +
  geom_point()
```

14.2

```
ggplot(a, aes(x = factor(a$site),a$orders)) +
  geom_point()
```

14.3

```
ggplot(a, aes(x = factor(a$platform),a$orders)) +  
  geom_point()+geom_smooth()
```

14.4

```
ggplot(a, aes(x = factor(a$platform),a$visits)) +  
  geom_point()
```