# Statistics Assignment

**-Meenu Jomi**

*Question 1:*

*The quality assurance checks on the previous batches of drugs found that — it is 4 times more likely that a drug is able to produce a satisfactory result than not.*

*Given a small sample of 10 drugs, you are required to find the theoretical probability that at most, 3 drugs are not able to do a satisfactory job.*

*a.) Propose the type of probability distribution that would accurately portray the above scenario, and list out the three conditions that this distribution follows.*

*b.) Calculate the required probability.*

Answer 1:

a.) The type of this distribution is known as **Binomial Distribution**.
   The three conditions for such a distribution are:
- The total number of trials are fixed.
- Each trial is binary, i.e. has only two possible outcomes, success and failure.
- The probability of success is the same for all the trials.

   They follow the following equation:
$$P(X=r) = {}^nC_r\,(p)^r\,(1-p)^{(n-r)}$$

b.)



(1-b)   $n = 10$

Using addition rule of probability:

$P(X \leq 3) = P(X=0) + P(X=1) + P(X=2) + P(X=3)$

$\boxed{P(X=r) = {}^nC_r\,(p)^r\,(1-p)^{n-r}}$

$P(X=0) = {}^{10}C_0\,(p)^0\,(1-p)^{+10}$

Finding value of $p$:-

| satisfactory | : | not satisfactory |
|:---:|:---:|:---:|
| 4 | : | 1 |

So, $P(\text{satisfactory}) = \dfrac{4}{5} = 0.8$

$P(\text{not satisfactory}) = 1 - 0.8 = 0.2$

$p = 0.2$

$$P(x=0) = {}^{10}C_0 \ (0.2)^0 \ (0.8)^{+10}$$

$$P(x=0) = \underline{0.107}$$

$$P(x=1) = {}^{10}C_1 \ (0.2)^1 \ (0.8)^9$$

$$P(x=1) = \underline{0.268}$$

$$P(x=2) = {}^{10}C_2 \ (0.2)^2 \ (0.8)^8$$

$$P(x=2) = \underline{0.302}$$

$$P(x=3) = {}^{10}C_3 \ (0.2)^3 \ (0.8)^7$$

$$P(x=3) = \underline{0.2013}$$

$$P(x \leq 3) = P(x=0) + P(x=1) + P(x=2)$$
$$+ \ P(x=3)$$
$$= \ 0.107 + 0.268 + 0.302 + 0.201$$
$$= \ 0.878$$

The required probablity that at most 3 drugs are not able to do a satisfactory job is around 87.8%.

_____

Question 2:

*For the effectiveness test, a sample of 100 drugs was taken. The mean time of effect was 207 seconds, with the standard deviation coming to 65 seconds. Using this information, you are required to estimate the range in which the population mean might lie — with a 95% confidence level.*

*a.) Discuss the main methodology using which you will approach this problem. State all the properties of the required method. Limit your answer to 150 words.*

*b.) Find the required range.*

Answer 2:

a.)  The main methodology used is known as **Central Limit Theorem**, also known as CLT.

This method is used when we try to find out a pattern or (in most of our case) probability of a large crowd (also known as **Population**) by randomly taking few from it (also known as **Sample**) and running all our test on the Sample to judge the Population's situation with minimal errors.

There are a few properties that help us to do the above stated task:

- Sampling distribution's Mean is equal to the population's mean.

Sampling distribution's mean = Population mean

$$\mu_{\bar{x}} = \mu$$

- Sampling distribution' Standard Deviation, also known as Standard Error is equal to the Population's Standard Deviation by the square root of the number of samples.

$$\text{Standard error} = \frac{\sigma}{\sqrt{n}}$$

$\sigma$ = Sample Population's Standard Deviation.

$n$ = Sample Number of samples taken.

= Sample size.

- If number of samples taken is more than 30 (n > 30), the sample distributions is a normal distribution.

b.)

2- b) Given,

No. of samples = $n$ = 100 nos.

Sample Mean = $\bar{x}$ = 207 sec.

Sample Sta Standard Deviation = $S$ = 65 sec.

Confidence level = y% = 95%

$$\text{Confidence interval} = \bar{x} \pm \left( \frac{Z^* \times S}{\sqrt{n}} \right)$$

From the $Z^*$ values table;
when confidence level $= y\% = 95\%$
$\Rightarrow$  $Z^* = 1.96$
$\therefore$ Confidence interval $= 207 \pm \left( \dfrac{1.96 \times 65}{\sqrt{100}} \right)$

i.e., Confidence interval $= (194.26 , 219.74)$

So, the range in which the population mean
lies is $(194.26 , 219.74)$.

(Z* value is also found by calculation.)

2-b) Finding of $Z^*$ by calculation:-

confidence level $= y\% = 95\%$

$\dfrac{y}{100} + \left( \dfrac{1 - \frac{y}{100}}{2} \right)$

$= \dfrac{95}{100} + \left( \dfrac{1 - \frac{95}{100}}{2} \right)$

$= 0.975$

From the z table we get

$Z^* = 1.96$

*Question 3:*

*a) The painkiller drug needs to have a time of effect of at most 200 seconds to be considered as having done a satisfactory job. Given the same sample data (size, mean, and standard deviation) of the previous question, test the claim that the newer batch produces a satisfactory result and passes the quality assurance test. Utilize 2 hypothesis testing methods to make your decision. Take the significance level at 5 %. Clearly specify the hypotheses, the calculated test statistics, and the final decision that should be made for each method.*

*b) You know that two types of errors can occur during hypothesis testing — namely Type-I and Type-II errors — whose probabilities are denoted by α and β respectively. For the current sample conditions (sample size, mean, and standard deviation), the value of α and β come out to be 0.05 and 0.45 respectively.*

*Now, a different sampling procedure (with different sample size, mean, and standard deviation) is proposed so that when the same hypothesis test is conducted, the values of α and β are controlled at 0.15 each. Explain under what conditions would either method be more preferred than the other, i.e. give an example of a situation where conducting a hypothesis test having α and β as 0.05 and 0.45 respectively would be preferred over having them both at 0.15. Similarly, give an example for the reverse scenario - a situation where conducting the hypothesis test with both α and β values fixed at 0.15 would be preferred over having them at 0.05 and 0.45 respectively. Also, provide suitable reasons for your choice (Assume that only the values of α and β as mentioned above are provided to you and no other information is available).*

Answer 3:

a.)

3. a)

$$H_0 : \mu \leq 200$$
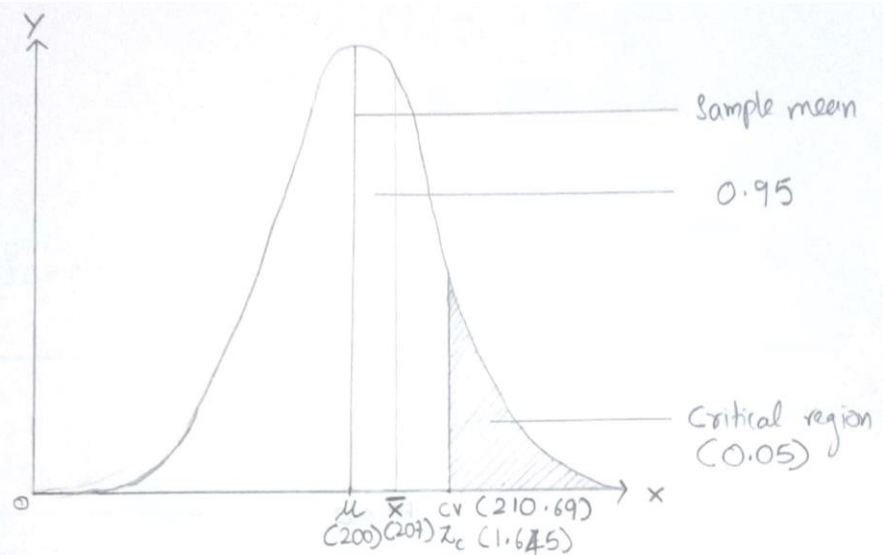$$H_1 : \mu > 200$$

(i) Critical Value Method :-

Standard $= \sigma = 65$ sec
Deviation

Sample $= n = 100$ nos.
Size

Signification $= \alpha = 0.05$
Level

Type of test :-

It is a one tail test, more specifically upper tail / right tail test.

Sample mean

0.95

Critical region
(0.05)

$\mu$ $\bar{x}$ CV (210.69)
(200)(207) $Z_c$ (1.645)

[ Graph for reference only. Not to scale ]

To find the value of $Z_c$ we need to find the probality at that point.

$$P = 1 - \alpha$$
$$= 1 - 0.05$$
$$P = 0.95$$

From the Z-table we find that, for P=0.95:

$$Z = 1.645$$

[ The value of 0.95 lies between the value of 1.64 and 1.65 so taking its average we get 1.645 ]

Since $Z = 1.645$

$\Rightarrow$ $Z_c = 1.645$

We will need the value of $\sigma_{\bar{x}}$ to find the critical value.

$$\text{Standard error} = \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

$$= \frac{65}{\sqrt{100}}$$

$$\sigma_{\bar{x}} = 6.\underline{5} = 6.5$$

Since it is a one-tail test we will only have one critical value (CV).

$$CV = \mu + (Z_c \times \sigma_{\bar{x}})$$

$$= 200 + (1.645 \times 6.5)$$

$$CV = 210.69 \text{ sec}$$

Given, $\bar{X} = 207 \text{ sec}$.
The value of the sample mean lies in the outside the critical region. So, we **fail to reject** the null hypothesis.

3- a)

(ii)    P-value method:

$$H_0 : \mu \leq 200$$
$$H_1 : \mu > 200$$

Sample Mean $= \bar{X}$ $= 207$ sec

Population Mean $= \mu$ $= 200$ sec

Sample Size $= n$ $= 100$ nos.

Standard deviation $= \sigma$ $= 65$ sec

Standard error $= \sigma_{\bar{x}} = \dfrac{\sigma}{\sqrt{n}}$

$$= \dfrac{65}{\sqrt{100}}$$

$$\sigma_{\bar{x}} = 6.5$$

Finding Z using the formula,

$$Z = \dfrac{\bar{X} - \mu}{\sigma_{\bar{x}}}$$

$$= \dfrac{207 - 200}{6.5}$$

$$Z = 1.077$$

To find the probability we need to use the
Z-table to find the values.

[Graph for reference only. Not to scale]



$$P(Z > 1.08) = 1 - P(Z < 1.08)$$
$$= 1 - 0.8599 \qquad [\text{from table}]$$
$$= 0.1401$$

So, p-value $= 0.1401$

Since it is one tail,

$$\text{Significance Level} = \alpha = 0.05$$

According to p-value method:

- If p-value less than $\alpha \Rightarrow$ reject null hypothesis
- If p-value more than $\alpha \Rightarrow$ fail to reject $H_0$.

So in this case,

$$p\text{-value} > \alpha$$
$$(0.1401 > 0.05)$$

Thus, we fail to reject the null hypothesis.

b.)

- A company creates pesticides and states that their pesticides are safe for environment. But lately we have had cases of newborns with deformations and the people believe it is due to these pesticides reaching the kids through their mothers while in womb.
  In this case the hypotheses will be:
  - $H_0$ : The pesticides does not contain any harmful chemicals.
  - $H_1$ : The pesticides contains harmful chemicals

  The Significance levels:

  - $\alpha$ should be very high ($\alpha = 0.95$)
  - $\beta$ will thus become very low ($\beta = 0.05$)

  <u>Reason:</u> This is because the issue is life threatening and needs to be assessed with maximum negligence.

  - If we follow the above levels, we can cause a Type 1 error. Which means the pesticides are safe and we state that it is unsafe.
  - If we do not follow the above levels, we can cause a Type 2 error. Which means the pesticides are unsafe and we state that it is safe.

  The first error is fine in this case. If we follow the second one, we are putting the lives of newborns at risk.

  We can't fix $\alpha$ and $\beta$ equal since we have a very serious life causing issue. So, it is best to have $\alpha$ high to avoid type 2 error.

- Deciding whether to take an umbrella while going to work, since the forecast states a small possibility of rain.
  In this case the hypotheses will be:
  - $H_0$ : It rains
  - $H_1$ : It does not rain

  The Significance levels:

  - $\alpha$ should be very low ($\alpha = 0.05$)
  - $\beta$ will thus become very high ($\beta = 0.95$)

  <u>Reason:</u> This is because getting wet in rain is not the right way to get till the office

  - If we follow the above levels, we can cause a Type 2 error. Which means we take an umbrella thinking it will rain and it does not rain.
  - If we do not follow the above levels, we can cause a Type 1 error. Which means we take an umbrella thinking it won't rain and it rains.

  The first error is fine in this case, since taking an umbrella and not raining won't affect anyone. But if we follow the second one there is a high chance of getting wet in the rain.

  We can't fix $\alpha$ and $\beta$ equal since doing so can cause us to get wet in the rain on the way to office. So, it is best to have $\alpha$ low to avoid type 1 error.
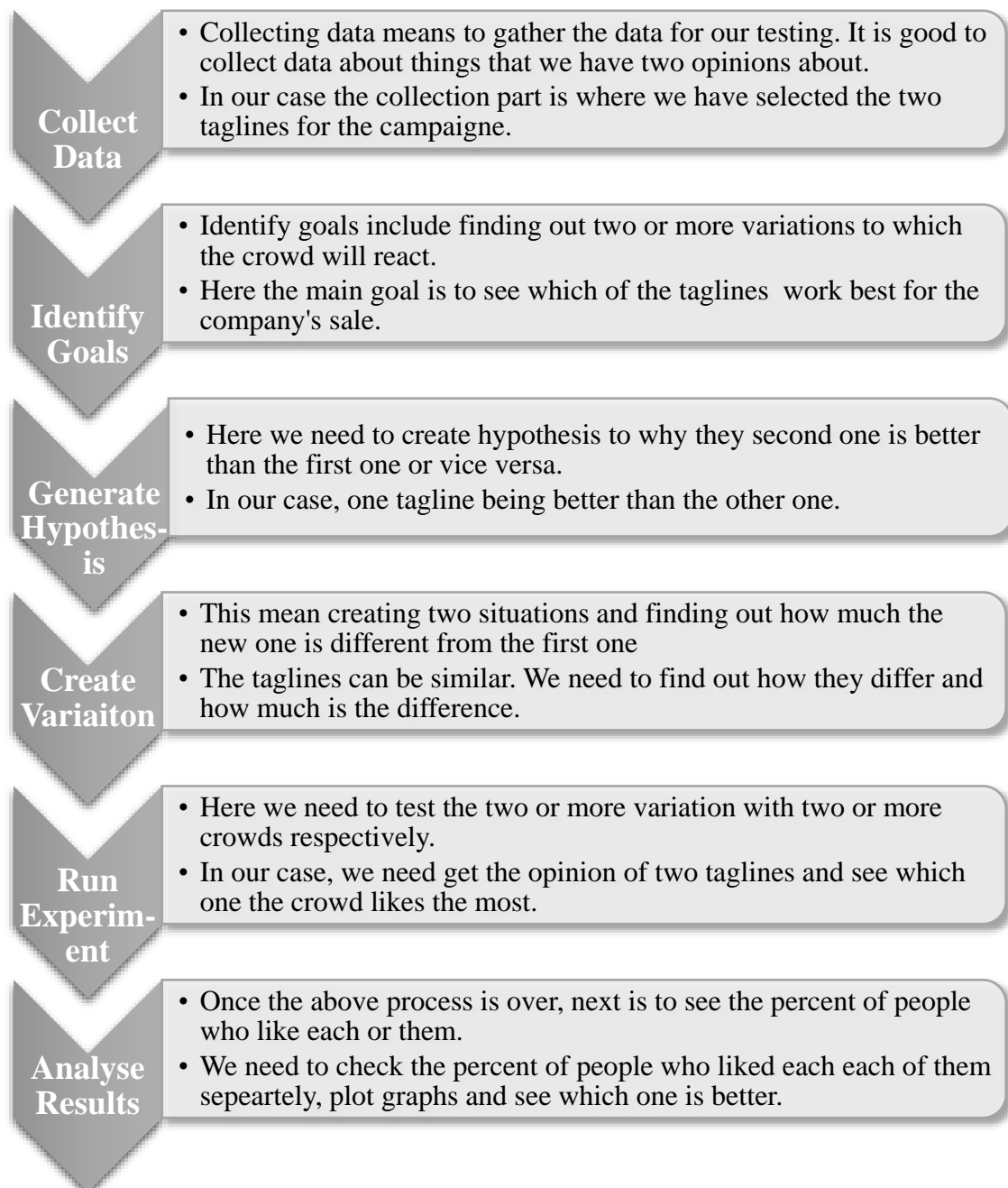
*Question 4:*

*Now, once the batch has passed all the quality tests and is ready to be launched in the market, the marketing team needs to plan an effective online ad campaign to attract new customers. Two taglines were proposed for the campaign, and the team is currently divided on which option to use.*

*Explain why and how A/B testing can be used to decide which option is more effective. Give a stepwise procedure for the test that needs to be conducted.*
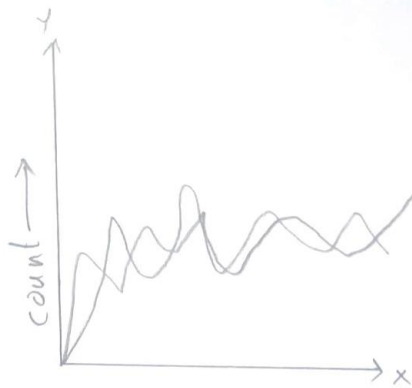
Answer 4:

An A/B testing have a few steps to follow. We sometimes repeat it again if we do not receive the expected results. In relation to our scenario we have do the following steps:

**Collect Data**
- Collecting data means to gather the data for our testing. It is good to collect data about things that we have two opinions about.
- In our case the collection part is where we have selected the two taglines for the campaigne.

**Identify Goals**
- Identify goals include finding out two or more variations to which the crowd will react.
- Here the main goal is to see which of the taglines work best for the company's sale.

**Generate Hypothesis**
- Here we need to create hypothesis to why they second one is better than the first one or vice versa.
- In our case, one tagline being better than the other one.

**Create Variaiton**
- This mean creating two situations and finding out how much the new one is different from the first one
- The taglines can be similar. We need to find out how they differ and how much is the difference.

**Run Experiment**
- Here we need to test the two or more variation with two or more crowds respectively.
- In our case, we need get the opinion of two taglines and see which one the crowd likes the most.

**Analyse Results**
- Once the above process is over, next is to see the percent of people who like each or them.
- We need to check the percent of people who liked each each of them sepeartely, plot graphs and see which one is better.
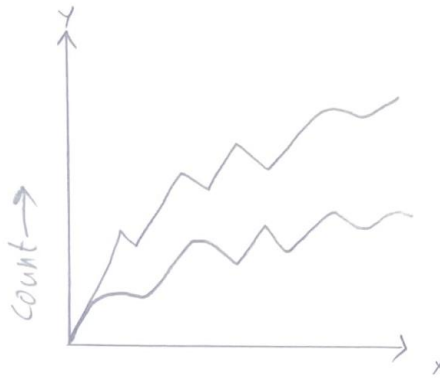
After the above we can get a graph:

- If the graph looks like below, we cannot come to a conclusion and might have to redo the entire process (It could also be because a very small sample is taken).

- But if this is the case, we can have a clear distinction and it is visible that that one above is preferred by most people.

Reason why A/B testing is useful in our case is,

1. once we have two variations, we can call them Version A and Version B,
2. put them both to test of a small crowd (called Sample) and see which of them works best for the public (called Population).
3. And from the results we will know for sure which of them to go for in the final presentation.

X------X------X------X