Test 2 covers the following topics:

The binomial distribution This describes the number of successes in n independent trials

The normal distribution The bell shaped curve

- The Central Limit Theorem This theorem describes the distribution of the sample mean  $\bar{x}$  as n, the sample size, gets large.
- **Sampling distributions** In general, a statistic summarizes a random sample an consequently is a random variable. A statistics distribution is known as a sampling distribution. We only spoke about  $\hat{p}$  and  $\bar{x}$ , but the topic is more general.
- Confidence intervals for p given  $\hat{p}$  Recall the  $(1 \alpha) \cdot 100\%$  confidence interval for p given  $\hat{p}$  is

$$\hat{p} - a\mathsf{SE}(\hat{p}) \le p \le \hat{p} + a\mathsf{SE}(\hat{p}),$$

where a is related to  $\alpha$  by the normal distribution:  $P(-a \le Z \le a) = 1 - \alpha$ , or  $P(Z \ge a) = \alpha/2$ .

Some sample problems I might ask would be: [Answers follow, using the R commands pnorm to replace the  $P(Z \le z)$  and dbinom(n,k) to compute P(X = k) for a binomial X]

1. If X is binomial with n = 5 and p = 1/4 find all of the following:

E(X), SD(X), P(X=3),  $P(X \le 1)$ 

```
Answer:

We have

> n = 5; p = 1/4

> EX = n * p; EX

[1] 1.25

> SDX = sqrt(n * p * (1 - p))

> SDX

[1] 0.9682

> PX3 = dbinom(3, n, p)

> PX3

[1] 0.08789
```

> PX1 = dbinom(1, n, p) + dbinom(0, n, p)
> PX1
[1] 0.6328

2. Suppose X is binomial with n = 4 and p = 1/2 and Y is binomial with n = 6 and p = 1/3. Which is more likely P(X = 2) or P(Y = 2)?

# Answer:

We have the same expectation (np = 2 for each, but for probabilities)

> dbinom(2, 4, 1/2)
[1] 0.375
> dbinom(2, 6, 1/3)
[1] 0.3292
Showing P(X = 2) is larger.

3. To see how effective text-messaging is for contacted students, 100 text messages were sent to 100 randomly chosen students. If the probability of being read is p = .75 compute the expected number read. Find the z score for 80 being read.

## Answer:

The expected number for a binomial is np = 75. The z-score is found by (x - mean)/sd or

- > (80 75)/sqrt(100 \* 0.75 \* (1 0.75))
- [1] 1.155
- 4. Let Z be a standard normal Find the following:

$$P(Z < 1), P(Z \le 2.3), P(Z \ge 1.23), P(-1 \le Z \le 1/2)$$

We use pnorm here instead of a table. You'd use a table on the test.

- > pnorm(1)
- [1] 0.8413
- > pnorm(2.3)
- [1] 0.9893
- > 1 pnorm(1.23)
- [1] 0.1093

> pnorm(1/2) - pnorm(-1)

- [1] 0.5328
- 5. Again, let Z be a standard normal. Find z for each

$$P(Z \le z) = .32, \quad P(Z \ge z) = 0.10$$

#### Answer:

These questions are about quantiles. We need to use qnorm() here as we don't have a table.

- > qnorm(0.32)
- [1] -0.4677
- > qnorm(1 0.1)
- [1] 1.282
- 6. Let Y be a normal random variable with mean 10 and standard deviation 20. Find

$$P(Y > 10), P(Y > 20), P(Y > 31), P(15 < Y < 25)$$

We find the *z*-score, then use pnorm().

- > 1 pnorm((10 10)/20)
  [1] 0.5
  > 1 pnorm((20 10)/20)
  [1] 0.3085
  > 1 pnorm((31 10)/20)
  [1] 0.1469
  > pnorm((25 10)/20) pnorm((15 10)/20)
  [1] 0.1747
- 7. Let  $X_1, X_2, \ldots, X_{16}$  is random sample for a normal population with mean 10 and standard deviation 20. Find the following

$$P(\bar{x} > 10), P(\bar{x} > 20), P(\bar{x} > 31), P(15 < \bar{x} < 25)$$

#### Answer:

Why is this different from the previous? Because the random variable is  $\bar{x}$ . The z scores need to be computed using  $\bar{x}$ 's mean and standard deviation. The mean is still  $\mu$ , but the standard deviation is now  $\sigma/\sqrt{25}$ .

[1] 0.1573

- 8. Suppose waist sizes are normally distributed with a mean of 92 cm and standard deviation of 11cm. Let Y denote a randomly chosen waist, find
  - (a)  $P(Y \ge 100)$ . (b)  $P(Y \ge y) = 0.80$

The first is done with a z-score

> 1 - pnorm((100 - 92)/11)

[1] 0.2335

The second is a quantile problem. First we find the z score of the quantile, then we convert into the right scale

> a = qnorm(1 - 0.8)
> 92 + a \* 11
[1] 82.74

- 9. Suppose  $X_1, X_2, \ldots, X_n$  is a random sample from a population with mean  $\mu$  and standard deviation  $\sigma$ . Which of these statements actually makes sense?
  - (a) The sample mean is the population mean.
  - (b) The mean of the sample mean is the population mean.
  - (c) the standard deviation of the sample mean is the population standard deviation.
  - (d) The distribution of the sample mean (for large n) is not the population distribution but the normal distribution.

### Answer:

The first is false – the sample mean is random, not a constant. The second is true. The third is false: it is  $\sigma/\sqrt{n}$ . Finally, the fourth is true by the Central Limit Theorem.

10. A survey of 365 Connecticut residents found a 60% supported current Senator Joe Lieberman. Find a 90% confidence interval for the population proportion.

We first find a using qnorm(). (Why do we use 0.95?)

> a = qnorm(0.95)
> a
[1] 1.645
> SE = sqrt(0.6 \* (1 - 0.6)/365)
> SE
[1] 0.02564
> MOE = a \* SE; MOE
[1] 0.04218
> c(0.6 - MOE, 0.6 + MOE)
[1] 0.5578 0.6422

11. A random sample of 50 people finds 46% agree with some proposition.

What is the margin of error?

Is .5 in the 95% CI based for  $\pi$  based on this sample?

How large would *n* need to be so that if  $\hat{p} = .46$  you would be sure that p = .5 is not in the 95% CI for *p* given by  $\hat{p}$ . (Set the margin of error to 0.04)

## Answer:

The margin of error depends on the confidence level. We'll use 95% so that we have

```
> SE = sqrt(0.46 * (1 - 0.46)/50)
> MOE = 1.96 * SE
> MOE
```

[1] 0.1381

Clearly 0.5 would be in a confidence interval. The margin of error informs us that 95% of the time  $\hat{p}$  within a margin of error of p.

If the MOE was less than 4 percentage points, then 0.5 would not be in the confidence interval. Using a formula from the notes, this will be the case is n is at least as large as

> a = 1.96; MOE = 0.04 > (a/MOE)^2 \* (1/4) [1] 600.2

- 12. The confidence interval for p based on  $\hat{p}$  requires a random sample from the population. Explain why this isn't the case in the following scenarios:
  - (a) The population is all US teens. The sample contains 1200 myspace users.
  - (b) The population is all people with colds. The sample contained only people who used homeopathic remedies.
  - (c) The population is all US teens. The sample is made up of the 120 friends listed by a myspace user.

# Answer:

Myspace users are not a random sample of US teens, presumably, so this population would introduce a bias into any calculation.

People who use homeopathic remedies do not necessarily represent the entire population of people who have colds.

A persons friends from myspace are in no way a random sample. An old adage "Birds of a feather flock together" might be interpreted in this case as saying the indivual sampling units are not independent.