

7th Class

6/15/10

“The race may not always go to the swift, nor the battle to the strong, but that’s the way to bet it.”--Damon Runyon

another take on correlation and causation [show cartoon from webpage]

last word on probability

[overhead on probability]

Jacob Bernoulli, posthumously published Ars Conjectandi in 1713; this is the beginning of the mathematical theory of probability--he developed the first law of large numbers

Ch. 6

random sampling

[overheads on sampling]

important insight: each individual observation in a random sample has the population probability distribution $p(x)$.

A sample is called a simple random sample (SRS) if each member of the population is equally likely to be chosen every time an observation is drawn.

A very simple random sample (VSRS) is a sample whose n observations X_1, X_2, \dots, X_n are independent. The distribution of each X is the population distribution $p(x)$; that is, $p(x_1) = p(x_2) = \dots = p(x_n) = p(x)$. Then each observation has the mean μ and standard deviation σ of the population.

This does not hold when the population is small and sampling is done without replacement. Note that nonreplacement is actually more efficient because each observation brings fresh information, but this gain in efficiency has to be weighed against the cost of having to use a more complicated procedure for estimating the mean

and standard deviation of the population, and in general is not worth doing/assuming unless the population is very small.

another important insight: Because of averaging, the sample mean \bar{X} is not as extreme (doesn't vary so widely) as the individuals in the population.

We can build up a sampling distribution for \bar{X} , denoted $p(\bar{x})$, by taking repeated sets of observations and calculating their mean.

note that $E(\bar{X}) = \mu$

this is because $\bar{X} = \frac{1}{n} [X_1 + X_2 + \dots + X_n]$

$$\text{so } E(\bar{X}) = \frac{1}{n} [E(X_1) + E(X_2) + \dots + E(X_n)]$$

$$= \frac{1}{n} [\mu + \mu + \dots + \mu]$$

$$= \frac{1}{n} [n\mu]$$

$$= \mu$$

also, $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$; since all the observations X_1, X_2, \dots, X_n are independent, we can use the formula from the last chapter to calculate the variance of a linear function of random variables (because of independence, no covariance terms to mess things up in extending this from 2 random variables to n random variables):

$$\bar{X} = \frac{1}{n} X_1 + \frac{1}{n} X_2 + \dots + \frac{1}{n} X_n$$

$$\text{so } \text{Var}(\bar{X}) = \left(\frac{1}{n}\right)^2 \text{Var}(X_1) + \left(\frac{1}{n}\right)^2 \text{Var}(X_2) + \dots + \left(\frac{1}{n}\right)^2 \text{Var}(X_n)$$

$$= \left(\frac{1}{n}\right)^2 \sigma^2 + \left(\frac{1}{n}\right)^2 \sigma^2 + \dots + \left(\frac{1}{n}\right)^2 \sigma^2$$

$$= \left(\frac{1}{n}\right)^2 (n\sigma^2)$$

$$= \frac{\sigma^2}{n}$$

and the standard deviation of $\bar{X} = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$

the standard deviation of \bar{X} is called the standard error of \bar{X} , or SE

note this goes to 0 as n goes to infinity

[note in the book they use the term moments of the sample mean]

Normal Approximation Rule, a.k.a. Central Limit Theorem

as n increases, the sampling distribution of \bar{X} becomes concentrated around μ and becomes more and more of a normal distribution, regardless of the underlying shape of the population distribution

[draw the distinction between:

the distribution of \bar{X} (with parameters $\mu, \frac{\sigma}{\sqrt{n}}$)

and the population distribution (with parameters μ, σ)]

1810 Pierre Simon Laplace came up with the Central Limit Theorem, where Central is meant as “fundamental”

repercussion of this theorem: can answer questions about how close the sample mean \bar{X} is to the population mean μ by using the standard normal table after transforming \bar{X} :

$$Z = \frac{\bar{X} - \mu}{SE}$$

Normal Approximation Rule for Proportions:

as n increases, the sampling distribution of P , the sample proportion, becomes concentrated around π and becomes more and more of a normal distribution, regardless of the underlying shape of the population distribution

the standard error of P is $\sqrt{\frac{\pi(1-\pi)}{n}}$, which goes to 0 as n gets large

so P is distributed with parameters $(\pi, \sqrt{\frac{\pi(1-\pi)}{n}})$

therefore, we can again use standard normal tables to describe how close the sample mean P is to the population proportion π by using the standard normal table after transforming P :

$$Z = \frac{P - \pi}{SE}$$

note that proportions are simply means for (0,1) (dummy) variables

do a couple of Ch. 6 type problems:

do problem 6-12:

mean 150, sd 25

capacity 7800

$n = 50$

a. what is the probability you end up overloaded?

can translate into a question about what is the average passenger weight

$$= \Pr(\bar{X} > 156) = \Pr(Z > (156-150)/(25/\text{SQRT}(50))) = \Pr(Z > 1.70) = .045$$

b. possible dependence in passenger weights (e.g. plane booked by a football team)

problem 6-36:

what is $\Pr(P > .50)$?

$$p_i = 43.2 / (34.9 + 43.2) = 0.55$$

$$SE(P) = \text{SQRT}[(.55)(.45)/1500] = .0128$$

$$\text{so } \Pr(P > .50) = \Pr(Z > (.50-.55)/.0128) = \Pr(Z > -.05/.0128) \\ = \Pr(Z > -4) = 1 - \Pr(Z > 4) \text{ (by symmetry)} = \text{basically zero}$$

QAC students can leave; a few review problems for Econ 300 students from old tests:

Person A tells the truth with probability $\frac{2}{3}$, person B with probability $\frac{3}{4}$, and person C with probability $\frac{4}{5}$.

Persons A and B make a statement which person C denies. What is the probability that the statement is true?

$$\Pr(\text{A tells the truth}) \cdot \Pr(\text{B tells the truth}) \cdot \Pr(\text{C lies}) = \frac{2}{3} \cdot \frac{3}{4} \cdot \frac{1}{5} = \frac{6}{60}$$

$$\Pr(\text{A lies}) \cdot \Pr(\text{B lies}) \cdot \Pr(\text{C tells the truth}) = \frac{1}{3} \cdot \frac{1}{4} \cdot \frac{4}{5} = \frac{4}{60}$$

$$\text{then } \Pr(\text{statement is true} \mid \text{A and B agree and C disagrees}) = \frac{\frac{6}{60}}{\frac{6}{60} + \frac{4}{60}} = 0.6, \text{ or } 60\%$$

Suppose that 3% of interviewees for a position with the National Security Agency lie on a certain question during a lie detector test. Suppose that 5% of the time the lie detector labels as a truth teller someone who is actually lying, and that 10% of the time the lie detector labels as a liar someone who is actually telling the truth.

- For a randomly picked person, what is the probability that the lie detector will label him or her a liar?
- If someone passes the lie test, what is the probability that the person is actually telling the truth?

This is an application of Bayes' Theorem:

	.90 "truth"	= .873
.97 truth	.10 "lie"	= .097
	.95 "lie"	= .0285
.03 lie	.05 "truth"	<u>= .0015</u>
		100% of possible states of the world

a. $\Pr(\text{"liar"}) = .097 + .0285 = 0.1255, \text{ or } 12.55\%$

b. $\Pr(\text{truth} \mid \text{"truth"}) = \frac{.873}{.873 + .0015} = 0.9982, \text{ or } 99.8\%$

Suppose that a drawing is to be done for four prizes. There are 100,000 ballots to choose from, but 30,000 of them contain a single person's name. What is the probability that the person will win:

- All four prizes?

- b. No prizes?
- c. At least one prize?

These can all be solved using the individual binomial probability table with $n = 4$ and $\pi = .30$

- a. $s = 4$, so $\text{Pr}(\text{win all prizes}) = 0.008$, or 0.8%
- b. $s = 0$, so $\text{Pr}(\text{no prizes}) = 0.24$, or 24%
- c. $\text{Pr}(\text{at least one prize}) = 1 - \text{Pr}(\text{no prizes}) = 1 - 0.24 = 0.76$, or 76%

A used car dealership has found that the length of time before a major repair is required for the cars it sells is normally distributed with a mean of 10 months and a standard deviation of 3 months.

- a. How many cars will still be running perfectly after a year?
- b. If the dealer wants no more than 10 percent of the cars to fail before the guarantee time, for how many months should the cars be guaranteed?

These are both solved using the standard normal probability table.

- a. $\text{Pr}(X > 12) = \text{Pr}(Z > \frac{12 - 10}{3}) = \text{Pr}(Z > 0.67) = 0.251$, or 25.1%
- b. $0.10 = \text{Pr}(Z > 1.28) = \text{Pr}(Z < -1.28) = \text{Pr}(\frac{X - 10}{3} < -1.28) = \text{Pr}(X < 6.16)$, so 6 months

At Daniel's Deli you can buy an 8 ounce or 12 ounce hamburger. You can also order a 16 or 24 ounce soft drink. Daniel has noticed that burger and soft drink combined orders have the following distribution:

	soft drink size	
hamburger size	16	24
8	.60	.15
12	.15	.10

- a. Calculate the mean and variance for hamburger size.
- b. Show that hamburger and soft drink sizes are dependent.
- c. What is the correlation between hamburger and soft drink sizes?

a. $E(\text{hamburger}) = 0.75 \cdot 8 + 0.25 \cdot 12 = \mathbf{9 \text{ ounces}}$

$\text{Var}(\text{hamburger}) = 0.75 + .25 \cdot 9 = \mathbf{3 \text{ ounces}^2}$

b. can show this several ways;

e.g., $\text{Pr}(\text{hamburger} = 8 \ \& \ \text{drink} = 16) = .6$

$\text{Pr}(\text{hamburger} = 8) \cdot \text{Pr}(\text{drink} = 16) = .75 \cdot .75 = .5625; .6 \neq .5625$

c. calculate $E(\text{drink}) = 18$; $\text{Var}(\text{drink}) = 12$, and Covariance (hamburger, drink) = 1.2

$$\text{so } \rho = \frac{1.2}{\sqrt{3 * 12}} = \mathbf{0.20}$$