# **Stat 100a: Introduction to Probability.**

Outline for the day:

- 1. Discuss midterm.
- 2.5.2 from hw2.
- 3. Hw3 notes.
- 4. Testing out your R function.
- 5. Conditional expectation.
- 6. LLN.
- 7. CLT.
- 8. CIs for  $\mu$ .
- 9. Sample size calculation.

Read through chapter 7. Remember hw3 is due Tue Sep 2 10 am. The project is due by email Tue Sep 2, 8pm.



1. Review of midterm problems.

#### 2. Problem 5.2 from hw2.

5.2. a) Y = the first time you are dealt a pocket pair or two black cards. Let p = P(get pair or 2 black cards) =  $13*C(4,2)/C(52,2) + C(26,2)/C(52,2) - 13/C(52,2) \sim 29.41\%$ . q = 1-p = 70.59%. Y is geometric with p = 29.41%, so P(Y=k) = q<sup>k-1</sup> p = .7059<sup>k-1</sup> \* 0.2941.

b) Z = the time til you've been dealt a pocket pair and you've also been dealt two black cards. P(Z > k) = P(X1 > k or X2 > k) = P(X1 > k) + P(X2 > k) - P(X1 > k and X2 > k). Let  $p_1 = P(\text{dealt a pocket pair}) = 13*C(4,2)/C(52,2) \sim 5.88\%$ .  $q_1 = 1-p_1 \sim 94.12\%$ .  $p_2 = P(\text{dealt two black cards}) = C(26,2)/C(52,2) \sim 24.5\%$ .  $q_2 = 1-p_2 \sim 75.5\%$ . P(Z > k) =  $q_1^k + q_2^k - q^k$ . P(Z > k-1) =  $q_1^{k-1} + q_2^{k-1} - q^{k-1}$ . P(Z = k) = P(Z > k-1) - P(Z > k) =  $q_1^{k-1} + q_2^{k-1} - q^{k-1} - [q_1^k + q_2^k - q^k]$ =  $q_1^{k-1} (1-q_1) + q_2^{k-1} (1-q_2) - q^{k-1} (1-q)$ =  $q_1^{k-1} (p_1) + q_2^{k-1} (p_2) - q^{k-1} (p)$ = .9412 <sup>k-1</sup> (.0588) + .755 <sup>k-1</sup> (.245) - .7059 <sup>k-1</sup> (.2941).

#### 3. HW3 notes.

Suppose X has the cdf  $F(c) = P(X \le c) = 1 - exp(-4c)$ , for  $c \ge 0$ , and F(c) = 0 for c < 0. Then to find the pdf, f(c), take the derivative of F(c). f(c) = F'(c) = 4exp(-4c), for  $c \ge 0$ , and f(c) = 0 for c < 0. Thus, X is exponential with  $\lambda = 4$ .

Now, what is E(X)?  $E(X) = \int_{\infty}^{\infty} c f(c) dc = \int_{0}^{\infty} c \{4exp(-4c)\} dc = \frac{1}{4}$ , after integrating by parts [ $\int u dv = uv - \int v du$ ] or just by remembering that if X is exponential with  $\lambda = 4$ , then  $E(X) = \frac{1}{4}$ .

For 6.12, the key things to remember are

1. f(c) = F'(c).

- 2. For an exponential random variable with mean  $\lambda$ ,  $F(c) = 1 \exp(-c/\lambda)$ .
- 3. For any Z,  $E(Z) = \int c f(c) dc$ . And,  $V(Z) = E(Z^2) [E(Z)]^2$ , where  $E(Z^2) = \int c^2 f(c) dc$ .

4. E(X) for exponential =  $1/\lambda$ .

5. E(X<sup>2</sup>) for exponential =  $2/\lambda^2$ .

Problem 7.14 refers to Theorem 7.6.8, p152.

You have k of the n chips in play. Each hand, you gain 1 with prob. p, or lose 1 with prob. q=1-p.

Suppose  $0 and <math>p \neq 0.5$ . Let r = q/p.

Then P(you win the tournament) =  $(1-r^k)/(1-r^n)$ .

The proof is by induction, and is similar to the proof of Theorem 7.6.6. Notice that if k = 0, then  $(1-r^k)/(1-r^n) = 0$ . If k = n, then  $(1-r^k)/(1-r^n) = 1$ .

For 7.14, the key things to remember are

- 1. If  $p \neq 0.5$ , then by Theorem 7.6.8, P(win tournament) =  $(1-r^k)/(1-r^n)$ , where r = q/p.
- 2. Let  $x = r^2$ . If  $-x^3 + 2x 1 = 0$ , that means  $(x-1)(-x^2 x + 1) = 0$ . There are 3 solutions to this.

One is x = 1. The others occur when  $x^2 + x - 1 = 0$ , so x = [-1 + / - sqrt (1+4)]/2 = -1.618 or 0.618.

So, x = -1.618, 0.618, or 1. Two of these possibilities can be ruled out. Remember that  $p \neq 0.5$ .

#### 4) TESTING OUT YOUR FUNCTION FOR THE PROJECT.

Suppose your function is called "neverfold". Run 6 neverfolds against 6 zeldas. install.packages(holdem) ## you must be connected to the internet for this to work. library(holdem)

a = neverfold

v = vera

decision = c(a,a,a,a,a,a,v,v,v,v,v,v)

name1 = c("n1","n2","n3","n4","n5","n6","v1","v2","v3","v4","v5","v6")

tourn1(name1, decision) ## Do this line a few times. Make sure there's no error.

### **5.** Conditional expectation, E(Y | X), ch. 7.2.

Suppose X and Y are discrete. Then E(Y | X=j) is defined as  $\sum_{k} k P(Y = k | X = j)$ , just as you'd think. E(Y | X) is a **random variable** such that E(Y | X) = E(Y | X=j) whenever X = j.

For example, let X = the # of spades in your hand, and Y = the # of clubs in your hand.a) What's E(Y)? b) What's E(Y|X)? c) What's P(E(Y|X) = 1/3)?

a. 
$$E(Y) = 0P(Y=0) + 1P(Y=1) + 2P(Y=2)$$
  
= 0 + 13x39/C(52,2) + 2 C(13,2)/C(52,2) = 0.5.

**b.** X is either 0, 1, or 2. If X = 0, then E(Y|X) = E(Y | X=0) and E(Y | X=0) = 0 P(Y=0 | X=0) + 1 P(Y=1 | X=0) + 2 P(Y=2 | X = 0) = 0 + 13x26/C(39,2) + 2 C(13,2) / C(39,2) = 2/3. E(Y | X=1) = 0 P(Y=0 | X=1) + 1 P(Y=1 | X=1) + 2 P(Y=2 | X = 1) = 0 + 13/39 + 2 (0) = 1/3. E(Y | X=2) = 0 P(Y=0 | X=2) + 1 P(Y=1 | X=2) + 2 P(Y=2 | X = 2) = 0 + 1 (0) + 2(0) = 0.So E(Y | X = 0) = 2/3, E(Y | X = 1) = 1/3, and E(Y | X = 2) = 0. That's what E(Y|X) is c. P(E(Y|X) = 1/3) is just  $P(X=1) = 13x39/C(52,2) \sim 38.24\%.$ 

# **6.** Law of Large Numbers (LLN) and the Fundamental Theorem of Poker, ch 7.3. David Sklansky, *The Theory of Poker*, 1987.

"Every time you play a hand differently from the way you would have played it if you could see all your opponents' cards, they gain; and every time you play your hand the same way you would have played it if you could see all their cards, they lose. Conversely, every time opponents play their hands differently from the way they would have if they could see all your cards, you gain; and every time they play their hands the same way they would have played if they could see all your cards, you lose."

Meaning?

LLN: If  $X_1, X_2$ , etc. are iid with expected value  $\mu$  and sd  $\sigma$ , then  $X_n ---> \mu$ .

Any short term good or bad luck will ultimately become *negligible* to the sample mean. However, this does not mean that good luck and bad luck will ultimately cancel out. See p132. 7. The Central Limit Theorem (CLT), ch 7.4.

Sample mean  $X_n = \sum X_i / n$ 

iid: independent and identically distributed. Suppose  $X_1, X_2$ , etc. are iid with expected value  $\mu$  and sd  $\sigma$ ,

$$\overline{X_n} \xrightarrow{\text{LAW OF LARGE NUMBERS (LLN)}}:$$

$$\overline{X_n} \xrightarrow{\text{CENTRAL LIMIT THEOREM (CLT)}}:$$

$$(\overline{X_n} - \mu) \xrightarrow{\div} (\sigma/\sqrt{n}) \xrightarrow{\text{CENTRAL Normal}}:$$

Useful for tracking results.



Truth: -49 to 51, exp. value  $\mu = 1.0$ 





day

Truth: uniform on -49 to 51.  $\mu = 1.0$ Estimated using  $\overline{X_n}$  +/- 1.96  $\sigma/\sqrt{n}$ = .95 +/- 0.28 in this example



<u>Central Limit Theorem (CLT)</u>: if  $X_1, X_2, ..., X_n$  are iid with mean  $\mu$  & SD  $\sigma$ , then  $(\overline{X_n} - \mu) \div (\sigma/\sqrt{n}) \longrightarrow$  Standard Normal. (mean 0, SD 1).

In other words,  $X_n$  has mean  $\mu$  and a standard deviation of  $\sigma \div \sqrt{n}$ .

Two interesting things about this:

(i) As  $n \rightarrow \infty$ ,  $X_n \rightarrow normal$ . Even if  $X_i$  are far from normal. e.g. average number of pairs per hand, out of n hands. X<sub>i</sub> are 0-1 (Bernoulli).  $\mu = p = P(pair) = 3/51 = 5.88\%$ .  $\sigma = \sqrt{(pq)} = \sqrt{(5.88\% \times 94.12\%)} = 23.525\%$ . (ii) We can use this to find a range where  $\overline{X_n}$  is likely to be. About 95% of the time, a std normal random variable is within -1.96 to +1.96. So 95% of the time,  $(\overline{X_n} - \mu) \div (\sigma/\sqrt{n})$  is within -1.96 to +1.96. So 95% of the time,  $(\overline{X_n} - \mu)$  is within -1.96  $(\sigma/\sqrt{n})$  to +1.96  $(\sigma/\sqrt{n})$ . So 95% of the time,  $\overline{X_n}$  is within  $\mu$  - 1.96 ( $\sigma/\sqrt{n}$ ) to  $\mu$  + 1.96 ( $\sigma/\sqrt{n}$ ). That is, 95% of the time,  $\overline{X_n}$  is in the interval  $\mu$  +/- 1.96 ( $\sigma/\sqrt{n}$ ).  $= 5.88\% + - 1.96(23.525\%)/\sqrt{n}$ . For n = 1000, this is 5.88% + - 1.458%. For n = 1,000,000 get 5.88% + -0.0461%.

# **Another CLT Example**

<u>Central Limit Theorem (CLT)</u>: if  $X_1, X_2, ..., X_n$  are iid with mean  $\mu$  & SD  $\sigma$ , then  $(\overline{X_n} - \mu) \div (\sigma/\sqrt{n}) \longrightarrow$  Standard Normal. (mean 0, SD 1). In other words,  $\overline{X_n}$  is like a draw from a normal distribution with mean  $\mu$  and standard deviation of  $\sigma \div \sqrt{n}$ .

That is, 95% of the time,  $\overline{X_n}$  is in the interval  $\mu$  +/- 1.96 ( $\sigma/\sqrt{n}$ ).

- Q. Suppose you average \$5 profit per hour, with a SD of \$60 per hour. If you play 1600 hours, let Y be your average profit over those 1600 hours. What is range where Y is 95% likely to fall?
- A. We want  $\mu$  +/- 1.96 ( $\sigma/\sqrt{n}$ ), where  $\mu$  = \$5,  $\sigma$  = \$60, and n=1600. So the answer is

 $5 + - 1.96 \times 60 / \sqrt{1600}$ 

= \$5 +/- \$2.94, or the range [\$2.06, \$7.94].

### 8. Confidence Intervals (CIs) for $\mu$ , ch 7.5.

<u>Central Limit Theorem (CLT):</u> if  $X_1, X_2, ..., X_n$  are iid with mean  $\mu$  SD  $\sigma$ , then  $(\overline{X_n} - \mu) \div (\sigma/\sqrt{n}) \longrightarrow$  Standard Normal. (mean 0, SD 1). So, 95% of the time,  $\overline{X_n}$  is in the interval  $\mu$  +/- 1.96  $(\sigma/\sqrt{n})$ .

Typically you know X<sub>n</sub> but not μ. Turning the blue statement above around a bit means that 95% of the time, μ is in the interval X<sub>n</sub> +/- 1.96 (σ/√n).
This range X<sub>n</sub>+/- 1.96 (σ/√n) is called a 95% confidence interval (CI) for μ.
[Usually you don't know σ and have to estimate it using the sample std deviation, s, of your data, and (X<sub>n</sub> - μ) ÷ (s/√n) has a t<sub>n-1</sub> distribution if the X<sub>i</sub> are normal.
For n>30, t<sub>n-1</sub> is so similar to normal though.]

1.96 ( $\sigma/\sqrt{n}$ ) is called the *margin of error*.



Over these 39,000 hands, Dwan profited \$2 million. \$51/hand. sd ~ \$10,000.

95% CI for  $\mu$  is \$51 +/- 1.96 (\$10,000 /  $\sqrt{39,000}$ ) = \$51 +/- \$99 = (-\$48, \$150).

Results are inconclusive, even after 39,000 hands!

## **9. Sample size calculation.** How many <u>more</u> hands are needed?

If Dwan keeps winning \$51/hand, then we want n so that the margin of error = \$51. 1.96  $(\sigma/\sqrt{n})$  = \$51 means 1.96 (\$10,000) /  $\sqrt{n}$  = \$51, so n = [(1.96)(\$10,000)/(\$51)]<sup>2</sup> ~ 148,000, so about 109,000 *more* hands.