Admin

- HW#4 almost done and LONG
most important the sea of ~~test~~ ~~test~~ genetics.

Plan

- a little bit more about CI's
- $\chi^2_{df}$ distr. (not on test)
- tests for variance
- ~~test~~ int of r-v.'s $(\hat{ch}.10)$
- cons. tables (ch. 5.1)(ch 5.2)
  (Cramis V not covered)

Sec 15.4  manipulating CI's

Let's say we have a CI for the profit made at one store and all stores
are essentially alike

$$CI_{M(\text{profit one store}), 95\%} = [\$8,000, \$9,000]$$

If we have 10 stores, we can just multiply by 10 to get a CI
for all stores

$$CI_{M(\text{all stores}), 95\%} = 10 \, CI_{M \text{ profit (one store)}, 95\%} = [\$80,000, \$90,000]$$

book talks about monotonic transforms — you'll see one on HW

What is the Chi-sq distribution?

$$\chi^2_{df} \stackrel{d}{=} Z_1^2 + Z_2^2 + \cdots + Z_{df}^2 \qquad \text{where } Z_1, Z_2, \dots Z_{df} \text{ independent}$$

$$= \left(\frac{X_1 - \mu_1}{\sigma_1}\right)^2 + \left(\frac{X_2 - \mu_2}{\sigma_2}\right)^2 + \cdots + \left(\frac{X_{df} - \mu_{df}}{\sigma_{df}}\right)^2$$

↑
Greek
letter chi, a squiggle

Turns out many things have the above configuration, so it is very useful.

Example: beer manufacturer making $0's. Need nozzle to be $\sigma = 0.2 oz$
and not higher otherwise there's overflow and underflow.
They run separate QA for the mean such as we did yesterday.
How to do this? Take an SRS $n=15$ beers at setup

test:

$H_0: \sigma = 0.2 oz \iff \sigma^2 = 0.04 oz^2$

$H_n: \sigma > 0.2 oz \iff \sigma^2 > 0.04 oz^2$

Take sample: $S = 0.27 oz \implies s^2 = .0729 oz^2$

Is this due to chance? Or is it too unlikely under the Null hyp?
Need test statistic with known distribution

$$Q \stackrel{\Delta}{=} (n-1)\frac{S^2}{\sigma^2} \sim \chi^2_{n-1} \qquad q = (25-1)\frac{.0729 oz^2}{.04 oz^2} = 43.74$$

so draw from $\chi^2_{24}$

proving this you need a PhD

If we look on sheet, at $\alpha = 0.05$, $\chi^2_{0.05, 24} = 36.4$

Since $q > 36.4 \Rightarrow$ reject $H_0$.

$pval = P\left(\chi^2_{24} \geq 43.74\right) = .0082 < \alpha = 5\% \Rightarrow$ reject $H_0$

Squarely is probably wrong with the machine

---

We're done with univariate statistics! that is, looking at data for one variable. e.g. all heights, all volumes, all wages for one company, etc

We now move on to bivariate analysis where we look at data between variables and how they relate. For instance:

How is a person's height associated with weight?

how the does gender affect SAT score $\longrightarrow$ interval vs. itval

$\longrightarrow$ nominal categorical vs. itval

how does Location Region affect purchase (ch 5) $\longrightarrow$

$\longrightarrow$ nominal categorical vs. nominal categorical

The rest of the class will be looking & analyzing bivariate data.

We start with categorical vs. categorical.

The way we show this data is through a "contingency table"

How

Purchase?

| | MSN | Reject Some | Yahoo | |
|---|---|---|---|---|
| Y | 6977 | 8282 | 5888 | 17,103 |
| N | 215 | 1 | 230 | 516 |
| | 7258 | 8283 | 6078 | 17,619 = n |

P. 78

"Cells" are mutually exclusive & collectively exhaustive of the data

leading to the "margins" or totals and represent the "marginal distr."

Why? Regardless of host, the # of people purchasing is:

| Y | 17103 |
|---|---|
| N | 516 |

The marginal distr. is the unconditional distr.

We can also see conditional distributions:

| | MSN | Reject Some | Yahoo | |
|---|---|---|---|---|
| Y | 96% | 99.98% | 96.22% | 97.07% |
| N | 3.93% | 0.08% | 3.26% | 2.93% |

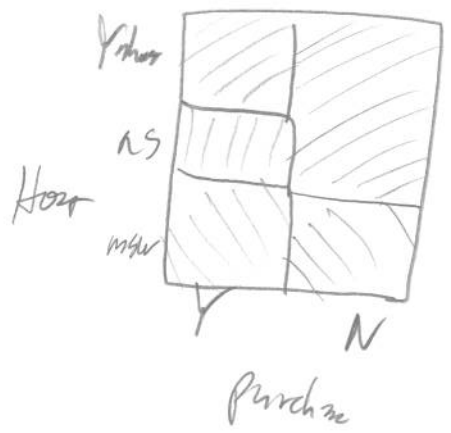| | MSN | RS | Yahoo |
|---|---|---|---|
| Y | 40.6% | 28.9% | 38.7% |
| N | 39.7% | 0% | 60.3% |
| | 81.2% | 28.8% | 34.5% |

Purchase | Yahoo

each column and row is a conditional distr.
You must calculate the % and figure difference in each case

It seems that the proportion of purchases varies based on the host. That means the host and the decision to purchase are <u>associated</u>.

There are different vocabulary words: correlated, dependent, causative. We'll get to those at the end of the week.

Another way to visualize data is with the segmented bar charts (p81) and mosaic plots (p82, 83). These charts make it easy to see association.
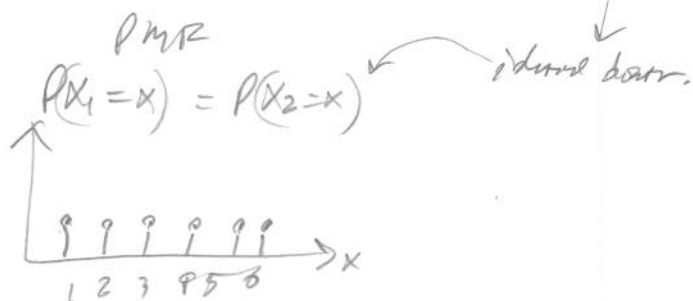


You can visualize associate much more easily

How strong is association?
We'd like to test the association if possible because there's always some due to chance...
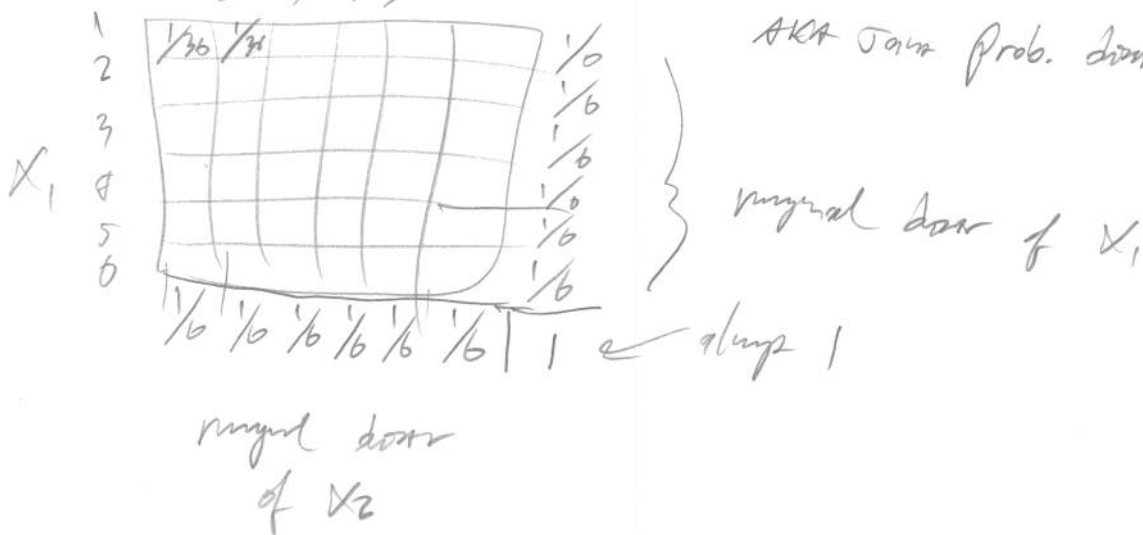
Let's review r.v.'s again... bit of ch.10.2

Two dice  $X_1, X_2 \overset{iid}{\sim}$  Dice

PMF

$P(X_1 = x) = P(X_2 = x)$ ← identical distr.



$$\begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \end{array}$$ x

What does independence really mean?

$X_2 \quad P(X_1 = x_1, X_2 = x_2)$ ← "Joint" PMF

AKA Joint Prob. distribution



marginal distr. of $X_1$

marginal distr. of $X_2$

← always 1

Independence is when $P(X_1 = x_1, X_2 = x_2) = P(X_1 = x_1) P(X_2 = x_2)$

for all possible states of $X_1, X_2$ ... (p.222)

Are our two dice independent above?

We essentially want to test whether or not variables are independent... If they're independent, the conditions must be the same as the marginals. For instance:   example in book p89-91

Poll 200 people   attitude towards shopping more online

$H_0$: Group, Attd $\perp$
$H_a$: ~
$\alpha = 5\%$

Observed

Attitude

|  |  | OK | Not OK |  |
|---|---|---|---|---|
| Group | Staff | 30 | 70 | 100 |
|  | Sales | 50 | 50 | 120 |
|  |  | 80 | 120 | 200 = n |

Expected

|  |  | OK | Not OK |  |  |
|---|---|---|---|---|---|
|  | Staff | 40 (20%) | 60 (30%) | 100 | (50%) |
|  | Staff | 40 (20%) | 60 (30%) | 100 | (50%) |
|  |  | 80 | 120 | 200 = n |  |
|  |  | (40%) | (60%) |  |  |

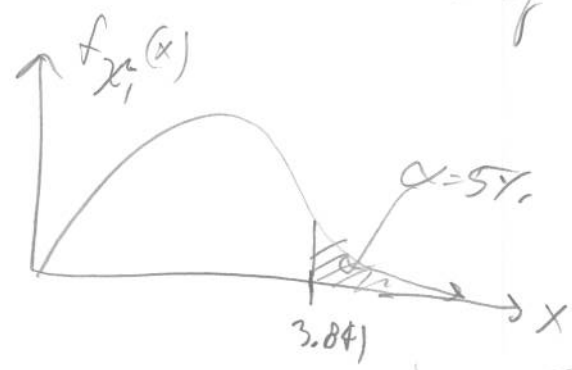non random   Under the null...
numbers

$$Q = \sum_{j=1}^{c} \sum_{i=1}^{r} \frac{(Obs_{ij} - Exp_{ij})^2}{Exp_{ij}} \sim \chi^2_{df} \quad \text{s.t. } df = (r-1)(c-1)$$
$$= (2-1)(2-1) = 1$$

$$q = \frac{(30-40)^2}{40} + \frac{(70-60)^2}{60} + \frac{(50-40)^2}{40} + \frac{(50-60)^2}{60} = 8.33 \quad \chi^2_{5\%,1} = > 3.841$$

which is a realization of $\chi^2_1$   $\Rightarrow$ reject $H_0$

$f_{\chi^2_1}(x)$

$\alpha = 5\%$

3.841   $\to x$

$pval = P(\chi^2_1 \geq 8.33) = .0039$
$\alpha = < 5\%$
$\Rightarrow$ reject $H_0$

# In-Class

Experiment: are two coin flips independent?

$H_0$: they are ind.

$H_a$: the are not ind.

$\alpha = 1\%$.

Observed        Coin 2

Coin 1



Expected        Coin 2

Coin 1



$$Q = \sum_{j=1}^{c} \sum_{i=1}^{r}$$     etc....

Bryce Cramer's V ... p96/93 (not cond)

Just because two variables are associated, does that mean
we can say there's causation? No

association $\nrightarrow$ Causation

Most subtle point in all of STATISTICS... probably
Most important message of semester... will do this tomorrow...