Lecture 24 (ch. 9, 11)

we leaved chi-squared test of k specific proportions in 1 pop.: How does proportion of X=01 $\#_{\mathcal{O}}: \mathcal{N}_{1} = \mathcal{N}_{\mathcal{O}_{1}}, \mathcal{N}_{2} = \mathcal{D}_{\mathcal{O}_{2}}, \cdots, \mathcal{N}_{n} = \mathcal{D}_{\mathcal{O}_{n}}$ something (e.g. tornadoes, X=1) Hi: At least 1 of These is wrong. vary across le categories chi-sqd dist. with df=k-1 (or k levels of a categorical var. Y)? 1 2-level X, 1 k-level Y. And The chi-squared test of independence: How does 1 categ. (X) He: X and Y are not indep. (affect another categ. / discrete var. (X) chi-sqd dist. with df= (k-1)(v-1)) XiY = 2 Cates. / discrite v.V.S. This test is equivalent to a "test of homogeneit" of population across categories " Shipped Shipped In regression we studied How does 1 (or more) continuous var (X) affect another continuous var. (Y) X,Y= 2 continuous vavs. How about How does 1 cates. (disevente var. (X) affect 1 continuous r.V. (Y)? This question requires comparing k means, ie. Mi=mean of Y for X=1, M2=mean of Y for X=2, ---- Mk=... X=k And The question of whether X affects Y becomes $H_{o}: \mathcal{M}_{i} = \mathcal{M}_{z} = --- = \mathcal{M}_{k} \qquad \left(N_{o}T \quad \mathcal{M}_{i} = \mathcal{M}_{ol}, \mathcal{M}_{z} = \mathcal{M}_{oz}, \cdots\right)$ H1: At least 2 m's are different. Once again, This method compares The mean of a continuous r.V. at different levels of a categ./Liscv. r.V. Example]: Does knowledge of religion depend on veligion?



Mormon		o J
hite evangelical Protestant	-0-	
White Catholic	-0-	March. eventhough we are testing
White mainline Protestant	-0	
Nothing in particular	—o —	the ther screval means are equal,
Black Protestant		in the state of the state
Hispanic Catholic 💙	nodel y (questag)	We must pay aftertion to variance i
	0 8 16	24 32
		Henrie The name ANOVA
		field he have his on 2

w

Example 2); Does fullnes of note shut have an effect on test scoves? mean test score note-sheet fullness not-so-full 0.6437 ١ 0.7205 Z 0.7179 3 0.7201 0.7142 VIVY full Again, looking at means is not enough. Must also look at variance. 35 30 0.8 22 (ouenbeu-20 0.6 tst2 12 8 0.4 2 Ó ŝ 0.2 0 2 з 4 5 2 1 з note-sheet fullness note-sheet 1.0 0.8 0.6 8t2 We are going to learn ANOVA F-test: F=0.5078p-v. 4.0 note that This test is just a generalization of The 2-sample/pop test (for comparing M, M2) to The case of k populations.



Does mean knowledge of religion vary across different religions? Does mean test grade vary across fullness of cheat sheet? Does mean vibration vary across different types of ball bearing? Does mean computer speed vary across different computers? Does mean detection error vary across different detection algorithms?

Generally: Are k means different? The Data will look like This 2 X (Catez) Note That we are now dealing with a generalization of the 2-sample t-test Ц 5 3 2 ۱ to more than 2 pops, Data l L ŧ here, 5 populations. 1 $H_0: M_1 = M_2 = \dots = M_5$ H1: Atleast Zu's are diff. The approach The Way ANOVA answers Cannot tell 5, that guestion is by finding out how much of The total variability in y is Within each categ./pop. Variabilit and how much is => M. +M27/43 within 5 between The categ. / pops. This "within", "butween" YI idea is a very powerful idea that shows-up Variability between every where!

total vaviability in the yij L'ategories/populations. L'ategories/populations. $n = \sum_{i=1}^{k} n_i$ - Grand mean $\downarrow \overset{k}{\underset{i=1}{\overset{n:}{\underset{j=1}{\overset{j}{1}{\overset{j}}{\overset{j}{1}{\overset{j}}{\overset{j}{1}$ $\sum_{i=1}^{k} \sum_{j=1}^{n_i} (\gamma_{ij} - \overline{\gamma})^2$ SST Z jth vergonse in The ith pop/ catez, k sample mean in ithpop. = (n:-1) sit in ith pop. h Ni $\sum_{i=1}^{k} N_i \left(\frac{\overline{Y}_i}{\overline{Y}_i} - \frac{\overline{z}}{\overline{Y}_i} \right)^2 + \sum_{i=1}^{k} \sum_{j=1}^{N_i} \left(\frac{Y_{ij}}{\overline{j}} - \overline{Y}_i \right)^2$ $\lim_{i \to 1} \frac{11}{55 \text{ between group}} SS_{within group} < 11$ SST = SSTr = Treatment ll SSE SSeephined SS un explained Variation within groups Variation between groups ~ Sample mean Jh Sample Variances ~ Sample var. of Sample means FYTIL + you see similarity with The ALOVA decomposition in vegression you are correct; There are similarities, and differences: \leftarrow of Herent les Ryresion df: n-1 = k + [n-(k+1)] \checkmark t h= # of B's

Now, we can compare SS between and SS vitain: Note: If all The $\frac{\text{Keovem}}{\text{If H}_{o}=\text{True}, F = \frac{SS_{bitween}/(k-i)}{SS_{within}/(n-k)} = \frac{MS_{bitween}}{MS_{within}}$ has an F-distribution with params. df = (k-l, n-k) All we need is Table <u>VIII</u> p-value = pr(F>Fobs) Again: The AMOVA Fitest is a generalization of The 2-sample t-test to move than 2 populations. One assumption of This Theorem is That The y's in each of The k populations are normal, and That They all have The same variance, ie. 012= 02= --= 02. (called homoscedasticity!) Use applots to test This [See optional har for] (understanding This

(example 9.) Consider 5 brands of computers. A code has been run on each of the brands 6 times, and The completion times have been recorded. Here are The data: Brand 1 2 4 3 13.1 15.0 14. 0 11. 6 $\overline{Y}_2 = 15.97$ 13.67 14.73 13.08 $S_2 = 1.167$ --- ---Ver. between = y = 13.68 W. within = 5, = 1.194 $\overline{\overline{Y}} = \sum_{i=1}^{5} \left(\frac{n_i}{n}\right) \overline{Y_i} = \frac{6}{30} (13.68) + \dots = 14.22$ Suturen = $\frac{5}{2}$ $N_i(\overline{Y_i} - \overline{\overline{Y}})^2 = 6(13.68 - 14.22)^2 + \dots = 30.88$ S_{ni} thin = $\sum_{i=1}^{n} \sum_{j=1}^{n} (Y_{ij} - \overline{Y_i})^2 = \sum_{i=1}^{n} (n_i - 1) S_i^2 = (6 - 1) (1 - 194)^2 + \dots = 22.83$ $F = \frac{30.88/(5-1)}{22.83/(30-5)}$ = 8.45 df=(5-1,30-5) See next page for how to read Table VIII. prule = pr(F7Fob) = pr(F)8.45) <.001 at 2=001 Conclusion: Rejut Ho (M=M2=...) in favor of H, (At least 2 u's are diff) which 2? In English: Brand has an effect on speed. Section 9.3 Shipped.

Reading Table VIII For the above example, we have dfnum=4, dfdenom=25 In Table VIII, for these df's you'll see Avea F 2.18 -> if Fobs = 2.18, Then p-value = 0.1 • 2.76 .05 4.18 }-> if Fobs = 3, Then no1 < p-value < .05 • 01 6.49 . 00 (] -> if Four=8, Then p-value 2,001 Most software produce an ANOVA Table (See pre lab) MS SS lf Source Fobs / P-Vila Between Group table VIII R-1 Ssbetween (factor) formla Wittin Gro-p n-k 85 within (ervor) SSTOTA n-1 Total For The above elample (from R): Source Sſ MS P-Value df F_{-} 30.85 7.71 8.44 5-1 81000. father 22.84 30-5 0.91 Evror 30 -1 53.7 Total

"large sample" "smill Z t smple" A Ox known/unknown Summary: Z,t H.: M=M. H.: M□M. 2 Ho: 7=70 H(: 7 D75 Z (No t). "large sample" $H_{o}: M_{1}-M_{2}=\Delta_{o} \quad H_{1}: M_{1}-M_{2} \quad \Box \Delta_{o}$ ₹,t indep. or paired Ho: $\pi_1 = \pi_2$, $\pi_2 = \pi_{32}$, $\dots \pi_n = \pi_{3k}$, H_1 : At least 1 is wrong Chised Ho: S2 Categ./ discrite Variables are independent Hi: not. Skipped. L2 pages are homogeneous wirt. le Categories chizd Ho: $M_1 = M_2 = \dots = M_k$ H: at least 2 m's are diff. F Again, note That The I-way AwovA problem deals with 1= {1 cont. var. (whose mean we are interested in), and x= [1 categ. var. (with le levels), > Nuct: Regression: 2 cont. variables, x, y.

(ch.11)

Wedid regression 1= x+Bx. + ··· + E; CR. 3. We did inference on M, 7, M,-M2, 7, -7, 7; ... (M=M2=M3=...) Ch. 7, 8, 9 Now, we do inference in regression (on B, d, Y(x), ...) Ch. 11. $\Rightarrow For a sample we write <math>Y_i = a + \beta x_i + \epsilon_i$ avbitvary params to be x_i estimated by OLS, ie. $a \neq SSE$, te. where $\hat{a}_{i}\hat{\beta}$ are The obs estimates of $\alpha_{i}\beta_{i}$, ie. $\beta = \frac{\overline{x_{Y}} - \overline{x_{\overline{Y}}}}{\overline{x_{\overline{Y}}} - \overline{x_{\overline{Y}}}} = \frac{S_{x_{\overline{Y}}}}{S_{x_{\overline{Y}}}}, \quad \widehat{\alpha} = \overline{Y} - \beta \overline{X}.$ where $S_{xx} = \sum_{i=1}^{n} (x_i - \overline{x})^2$ Recall That $S_{xy} = \sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})$ Sample $V_{av} = S_x^2 = \sum_{n=1}^{n} (x_i - \overline{x})^2 = \frac{S_{xx}}{n-1}$ > The Analysis of Variance: SST = Z(7:-7) = SSeeplaind + SSurexpl. SSE $df_{1} n-l l = k + (n-(k+l))$ R²= SST # of p's $S_e = \sqrt{\frac{SSE}{n - (k+1)}} \sim RMSE$ (chelinding d) percent of variability in y explained by x --std. dev. of errors to predictors ~ Typical evvor or spread about fit. Y= d+ B, X, + ... Buxk (Goodners of fit)

> (For population

Now, let's consider the population. For the moment suppose we have it, Just because we have The population, it does not follow That There is no scatter between a, y. I.e. even for The population, There is a scatter between x and y, and so There is an ous fit for the pop! > Ic. even for the population there is an OLS fit! Call it the true fit. what symbols should we use to denote this true fit!

Sample mean : X Sample OLS fit : X, B true/pop/distr. mean: Mx true/pop/distr. OLS fit: ?

It's tempting to use a, B (w/o hat). But a, B are supposed to be free parameters that are tuned to minimize The SSE. So, technically, we should introduce new symbols for The true OLS fit, However, to keep Things simple, we will go ahead and use a, B to denote The true ous fit. I.e. up to now, a's have been free pavameters to do da, dip, ... But hence for the dis denote The OLS fits for the population.

= In short: $\hat{\gamma}(x) = \hat{\alpha} + \hat{\beta} \times (for sample)$ and $\gamma(x) = \alpha + \beta \times (for population)$

will be used for The respective predicted values. I.e. ŷ(x) = prediction (in sample)

y(x) = prediction (in pop.)

 \Rightarrow

Now, to do inference we need a probability model (for regression): Assume is are Normally distr. at each x, with params M= Y(x), O= OE e.g. $M = \Psi(x) = a + \beta x + \dots$ $U = U_{\mathcal{E}} = fixed$ estimate a, β with $\hat{a}, \hat{\beta}$ estimate with $\frac{s}{2}$ $\frac{\gamma_{1}}{\gamma_{1}}$ FTI At a given x, YNN(Y(x), JE) $\langle \varphi \rangle = \gamma(x) + E \implies E = \gamma - \gamma(x) \sim N(0, \sigma_E)$ In short: For a This allows us to say things like: fixed x, everything 1) <u>At a given x</u>, $\hat{Y}(x)$ estimates the true mean of Y, Y(x) we have done (CI, $\hat{x}_{+}\hat{\beta}x$ applies to Y. p-value...) non applies to y. 2) At a given x, we expect about 25% of The y's to be within $\chi(x) \pm 1.96 \sigma_E$. like 95% of x's are within 14 ± 1.965 (Ch.1) Y(x)-1960 Y(x) Y(x)+1960 not x because. 3) At a given x, we can find other proles for y. e-j- prob(a<y<b (x) = fike pv(a<x<b) = (Ch.1)prob (a-y(x) < (y-y(x)) / D- y(x)) = Table T $\overline{\sigma_{e}}$ $\overline{\sigma_{e}$ $P\left(\frac{\alpha-\mu}{\sigma} < \frac{x-\mu}{\sigma} < \frac{b-\mu}{\sigma}\right)$ Note that these are just prob calculations, not p-values or CIS.

hw_lect24_1: The following data refer to the melting temperature, y (in some unit), of a certain material at four different pressures, x (in some unit).

P	ressu	are Temperature					
-							
1	.6	59.5.53.3.56.8.63.1.58.7					
3	.8	55.2, 59.1, 52.8, 54.4					
6	.0	51.7, 48.4, 53.9, 49.0					
1	0.2	44.6, 48.5, 41.0, 47.3, 46.1					

a) Make a comparative boxplot of y for the four pressure levels. By R.

b) Based on the above boxplot, would you say there is a difference in the mean melting temperature for at least 2 of the pressure levels?

c) At alpha = 0.05, is there evidence that the mean melting temperature at the at least 2 of the four pressure levels are different? Report the p-value, and state the conclusion clearly. By R; see prelab to see how to do 1-way ANOVA in R.

d) Write code to compute the above p-value "by hand," i.e. without using aov() or lm(), but using the basic formulas for SS_between, SS_within, etc.

e) After (or before) a 1-way ANOVA test, one should check the two assumptions that the y's are normally distributed within each group, and with the same variance. To that end, make a plot that shows four qqplots (one for each pressure level) superimposed onto a single figure; make sure that the four qqplots have different colors. Are the 4 qqplots reasonably straight, and do they have approximately equal slopes? Hint: in the first call to plot(), use xlim=c(-2,2) and ylim=range(y). Use the "By hand" method for making qqplots.

hw_lect24_2

In a problem dealing with flow rate (y) and pressure-drop (x) across filters, it is known that y = -0.12 + 0.095 x. Note: this is the true "fit" to the population. Suppose it is also known that sigma_epsilon = 0.025. Now, IF we were to make repeated observations of y when x=10, what's the prob. of a flow rate exceeding 0.835? Y = -0.12 + 0.095 x, $V_E = 0.025$

hw_lect_optional:)By hand or by R (see prelab)

Consider the data you collected. Take one of the continuous variables (call it y) and the categorical (or discrete) variable with 3 or more levels (call it x). Since x is discrete/categorical, we can consider each level as a different population. Eg, if your x has 3 levels (say H, M, L), separate the corresponding y's into 3 populations.

a) Do 1-way ANOVA to test if any of the k populations have different means. Report the p-value, and the conclusion In English.

b) Make qq-plots of y for each of the levels of X. It's important to have all qq-plots superimposed onto a single figure. See "By hand" q-plots in prelab. For example, if your x has 3 levels, then you need to have 3 qq-plots superimposed on one plot, e.g.

 \mathbf{i}

Recall that equal-slopes translate to equal variances, and so this will be a way of visually checking the homoscedasticity assumption mentioned above (in the lecture note).