Lecture 6 (Ch. 1, 2) גלסיבן Last time we learned how to find areas under N(0,1). left-areas What about P(K, 0)? TableI It would be impractical to have tables for all 11 and o! Fortunately, there is a trick: Change variables! also called standardization X-> Z = X-M (Z-SCOVER) Normal with Pavams M. 0 A-M Normal with A-M Normal with Normal wi So, to compute area between 2 values: Important first step: $prob(a < x < b) = prob(\frac{a-\mu}{\sigma} < \frac{x-\mu}{\sigma} < \frac{b-\mu}{\sigma})$ Standardize! graphically = prob (Z < b-m) = prob (Z < a-m) Algebraical Eitherway (algebraically or graphically) you can obtain The value of each term from Table 1. (Example: If x~ N(M=4, J=3), What's pr(-2<x<2)? $\frac{\sigma}{2} = \frac{\sigma}{3} = \frac{\sigma}{2} = \frac{\sigma}$

Standardizing with $x \rightarrow z = x - m$ is specific to $N(\mu, \sigma)$: $\chi \sim \mathcal{N}(\mu_{1}\sigma) \longrightarrow Z \sim \mathcal{N}(0,1)$ The main point is that we can then find pr((x(): (See hur) $pr(a < x < b) = pr(\frac{a-\mu}{\sigma} < 2 < \frac{b-\mu}{\sigma})$ But, more generally, any change of variable, x->y, That allows you to compute pr(--<x<--) is useful (and can be called "standarditing").

distribution percentile (or quantile) So for: Find area, given x Now : Find x, given (some left/right) avea. The simplest example is The distribution median: area= 0.5 = 50% E.g. metion: $\int f(x) dx = \frac{50}{100}$ Median is a special case of the more general concept of percentile: $f(x) dx = \frac{n}{100}$ f(x) 1 nth percentile median = 50 percentile = 0.5 quartile = 2nd quartile Example;) what's The 90th percentile (0.9 quantile) of N(M, J)? Graphically: 90% ? = M + 1.2850 $\in 1.285 = \frac{?-M}{2}$ Table IAlgebraically: $0.9 = pv(x < ?) \stackrel{=}{=} pv(\frac{x-\mu}{\sigma} < \frac{?-\mu}{\sigma}) = pv(\frac{z}{z} < \frac{?-\mu}{\sigma})$ Table I $\rightarrow \frac{2-1}{2} = 1.285$?= 1+ 1.2850

Important comments:

1) The nth percentile/quantile/quartile/... is a number on the x-axis in the above figure. In other words, it has the same units as x itself. If x is in Kg, then the nth percentile is a number in Kg. It is NOT (necessarily) a percent.

2) The notion of the nth percentile/quantile/quartile/... applies to histograms, as well. In that case, it's called the sample percentile/quantile/quartile/... . So, for example, the 65th sample percentile of data on weight is a specific value of weight in your sample for which 65% of the cases are smaller.
3) The notion of the nth percentile/quantile/quartile is an extremely useful concept. It shows-up everywhere. One place is in summarizing distributions and histograms. Look:

Suppose you want to find out which of two computers is faster. you take a given program, and run it on each computer 100 times, and record the times it takes to run the code to completion. You can/should then look at the histogram of "completion time" for the 2 computers!

computer A computer B- \neq completion time = χ

Which computer is faster? Discussion:

1) It looks like B is faster on the average. But B is also more moody! (Learn to also look at the width of histograms.) So, which computer should you buy?! The answer: It depends.

2) More importantly, the huge overlap between the histograms causes a huge problem. The true mean of x for the two computers is somewhere close to the center of each histogram; but we don't know where After all, our data is only a sample taken from some unseen population. So, if/when there is too much overlap, then we cannot tell which computer is truly faster; see the next page, too.

Comparison of even 2 groups is complex, involving location, width, ... How do we handle more Than 2 groups? Quartiles are The basis of The so-called 5-number summary of a hist (or dist); often plotted as a boxplot: This monterio histogram -.25 from 2.3 b .25 fite better have 2ⁿ¹quartile for let & Lob



Observations: Based on this sample, we cannot really tell if one computer is faster. We could have been able to say something if one hist/boxplot were strictly shifted with respect to the other. Note the the REASON we caanot tell is the WIDTH of data/hist/boxplot. This is one way in which the width of data plays an important role.

Now. suppose we decide to compare the computers only in terms of the center (say, median or mean). Then, computer B is faster "on average" because its typical completion time is shorter. But computer B is also more "moody" (less consistent), because it has a wider spread in completion times. So, again, WIDTH plays an important role.

Having said all that, one cannot conclude that computer B is faster, because these boxplots are based on a sample/ hist. We do not know what is the distribution of x (i.e. the population). We do know the true dist./population mean (or median) of x for each computer is somewhere in the boxplot, but we don't know where. Given the huge overlap between the boxplots, we cannot conclude that B is faster in the population. In fact, in this case we cannot conclude anything about the population/distribution because there is too much overlap. Get used to this kind of conclusion! It happens because of the existence of WIDTH!

How much overlap is too much? Ans. in Ch-7 and beyond. For now, just learn that every time you see a number, it's actually a sample (of size 1), and that it's actually a single realization of a random variable, and that the variable actually has a spread/width.

Here is an example involving many groups probably Cannot Till d'ffere 21011 0.8 test . 0.5 Do not Think That as a math major you are destined to do poorly in 390 0.4 This information should be used to improve Things, 0.3 0-C-EPRMJ 0-C-HIST 0-C-PHIL 0-C-PHIL 0-C-PHVS 0-C-PREMAJ 0-D-ATM-S 0-D-ENVIR 0-E-MKTG 0-J-A-A 0-C-BIOCHM 0-C-C-SCI 0-A-EEP -A-N-MATR -A-UWACAD 0-C-ACMS 0-C-CHEM 0-C-PRESCI D-J-EXPENG 0-J-CMP-E 0-J-E-E 0-J-M-E -J-PRENGF D-O-BIOEN class discussion Note That The discussion involves comparing The whole box plots, not comparison of The 5 numbers one by one In summary, (comparative) box plots form a powerful tool of Viscolly comparing multiple groups in terms of either data/sample from each group or Their distributions.

hw-lette-1) If x follows N(M, o), what's The prob. of x being within 1 or of M? Get used to This jargon. the-lect 6-2) Standardization is important in finding probs. Although it almost always refers to The change of variable $z = \frac{x-m}{5}$, taking $N(\mu, \sigma)$ to N(0, 1), sometimes a different change of variable is required to obtain something That has N(0,1) distr. Find The pr(x(2)) if $(\frac{1}{x})$ has a std. Normal dist. the let 6-3) Another of The named distributions is The so-called power-law dist. It's formula is f(x) = dxd-1, o(x<1, where 20 is its pavameter. a) Find The nth percentile of that distribution. Hint: The answer will depend on 'a and n. b) How long is The box portion of The corresponding boxplot? Hint: The answer will depend on a.

hw lect6-4

Consider ONE of the 2 continuous random vars, and ONE of the 2 discrete/categorical variables in the data you collected. Make comparative boxplots for the continuous variable for each level of the discrete variable. E.g. if the discrete var. has 4 levels, then you need to show 4 boxplots for the continuous variable, all on the same plot, side-by-side. Interpret/discuss them. By R.