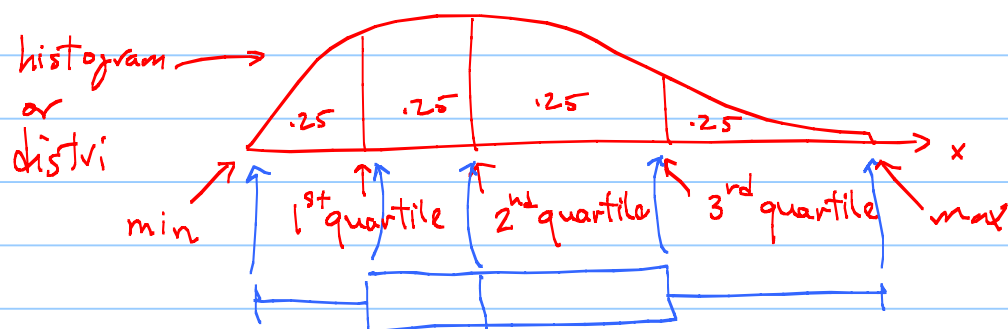


Lecture 6 (Ch. 1 mostly)

Last time we introduced the concept of the n^{th} percentile (for the normal distr.); an x value with $n\%$ area to its left. Note that percentiles (or quantiles, quantiles, ...) apply to distrs and hists.

Quartiles are the basis of the so-called "5-number summary" of a hist (or distr), often plotted as a boxplot:

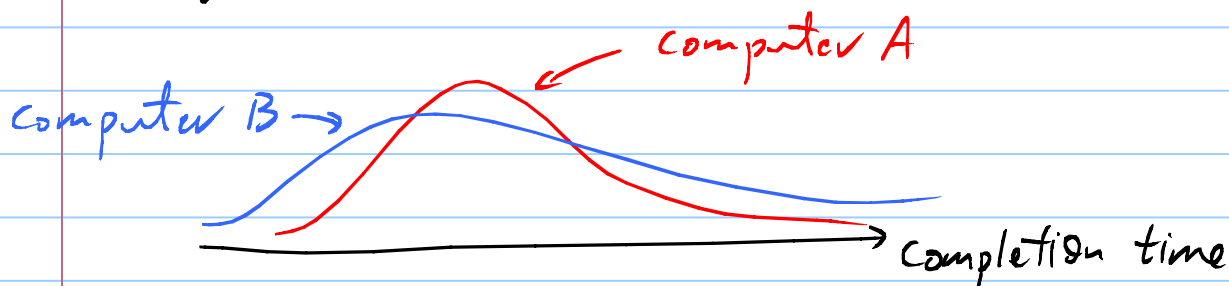


This material is from 2.3, but fits better here for Lect. & Lab.

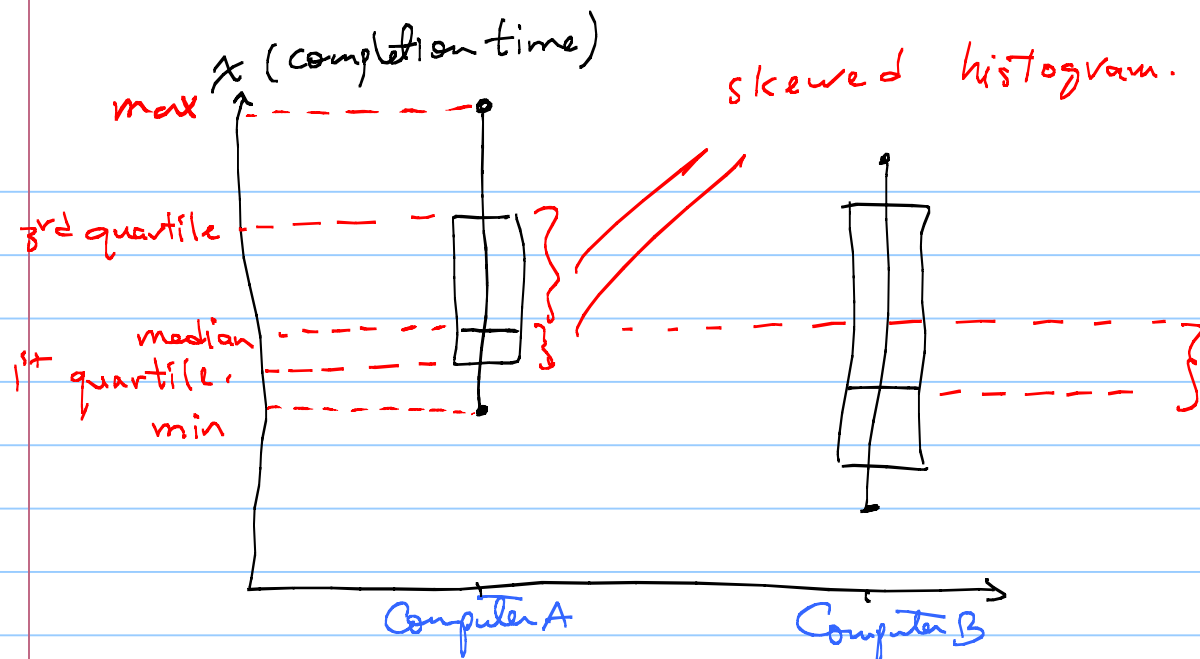
2nd quartile = 50th percentile = median = splits data in half.

1st (3rd) quartile = median of 1st (2nd) half.

Eg. 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 then look at the histogram of "completion time" for the 2 computers:



The interpretation of such results is complex (see next page). Boxplots allow us to handle problems like this even involving many more (than 2 or 3) computers.

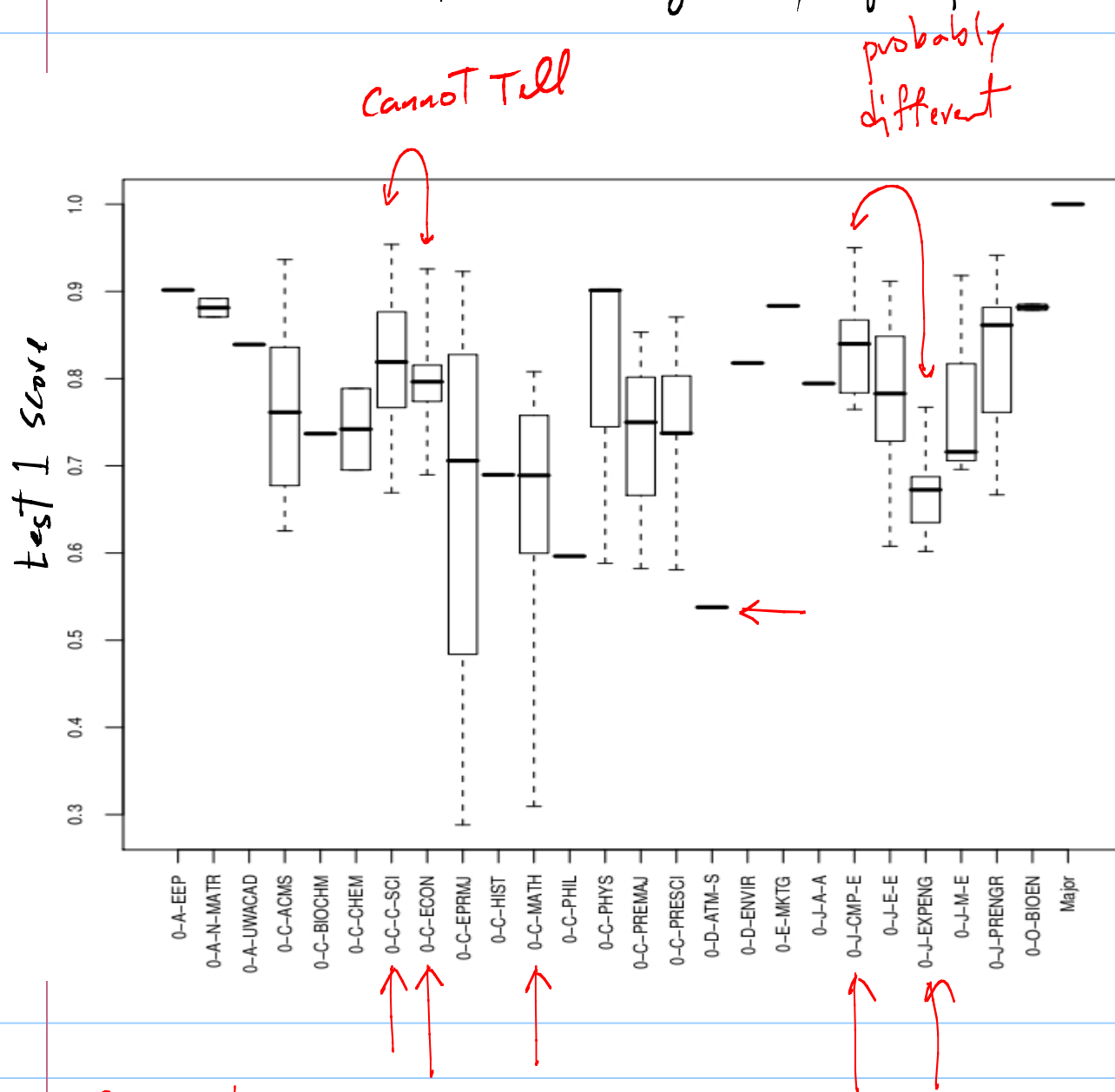


Observations: Based on this sample, computer B is faster "on average" because its median completion time is shorter. But computer B is also more "moody" (less consistent), because it has a wider spread in completion times. Important: Note spread!!

Having said all that, one cannot conclude that computer B is faster, because these boxplots are based on a sample of size 100. We do not know what the true distribution of x is. The True / 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. We cannot conclude anything! How much overlap is too much? Ans. in Ch. 7, 8.

For now, just learn that everytime 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. And that's important!

Here is an example involving many groups :



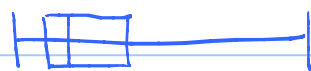
class discussion!

Note that the discussion involves comparing the whole boxplots, not comparison of the 5 numbers one by one.

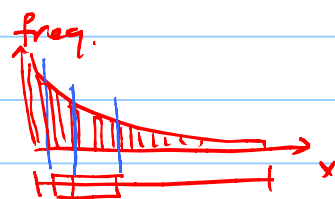
In summary, (comparative) box plots form a powerful tool of visually comparing multiple groups in terms of either data/sample from each group or their distributions.

Here is a clicker question from the post:

Consider this boxplot of a histogram
The histogram itself will "look"



A) Uniform B) Normal C) Exponential D) cannot tell



In the above question, one can also conclude that the population (i.e. distribution) from which the sample was drawn may be exponential.

Recall that we use dists. to represent populations, and hists to represent the sample/data from that pop.

We have been using dists. as mathematical objects.
And they are! But it may help to derive one.
Next.

Derivation of Binomial:

Consider N objects (population), where

Each object is 1 (Head, Girl, ...) or 0 (Tail, Boy, ...)

Suppose the proportion of 1's in the pop. is known = π .

Now, select n (e.g. 3) of the objects (with replacement) = sample and note the value of each object.

Repeat many many times (e.g. 10^8)

Q What proportion (of the 10^8) will be 1,1,1? 1,1,0? Etc.

Note: I'm not asking for the prop. of 1's in each sample.

I'm asking for the prop., out of the 10^8 trials, that are 1,1,1. Etc. ie. large!

A

		$X = \# \text{ of } 1\text{'s}$
prop. of 1,1,1	$\pi \cdot \pi \cdot \pi$	3
1,1,0	$\pi \cdot \pi \cdot (1-\pi)$	2
1,0,1	$\pi (1-\pi) \pi$	2
0,1,1	$(1-\pi) \pi \pi$	2
Etc.		
0,0,0	$(1-\pi) (1-\pi) (1-\pi)$	0

independence

prop($X=3$) = $1 \pi^3$

prop($X=2$) = $3 \pi^2 (1-\pi)$

prop($X=1$) = $3 (1-\pi)^2 \pi$

prop($X=0$) = $1 (1-\pi)^3$

$\frac{3!}{3! (3-3)!}$
 $\frac{3!}{2! (3-2)!}$
 $\frac{3!}{1! (3-1)!}$
 $\frac{3!}{0! (3-0)!}$

$\therefore \text{prop}(X=x) = \frac{3!}{x! (3-x)!} \pi^x (1-\pi)^{3-x}$

$x = 0, 1, 2, 3$

$$\therefore \text{prop}(X=x) = \frac{n!}{x!(n-x)!} \pi^x (1-\pi)^{n-x} \quad (\text{Table II})$$

$x = 0, 1, 2, \dots, n = \# \text{ of } 1\text{'s out of } n$

This is the mass function, $p(x)$, of a binomial variable x .
 E.g. $x = \# \text{ of heads out of } n \text{ tosses}$

Because we derived the above expression using proportions, it follows that $\sum_x p(x) = \sum_x \text{prop}(x) = 1$.

Recall The connection between coin tosses and sampling:

The prob. of getting x heads out of n tosses of a coin
 (or 1 toss of n coins)

The prob. of getting x boys out of a sample of size n .
 " " " x defective gates on a chip with n gates
 Etc

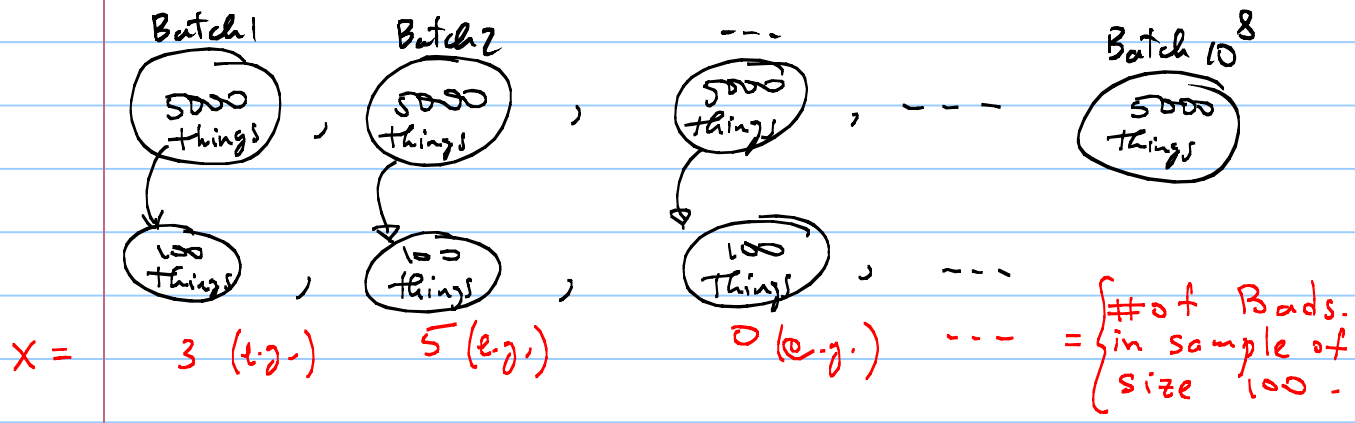
What's π ?

For the coin example, it's the prob. of getting a H on one toss.
 In the other example, it's the prob of drawing a boy,
 i.e. the proportion of boys in the pop.

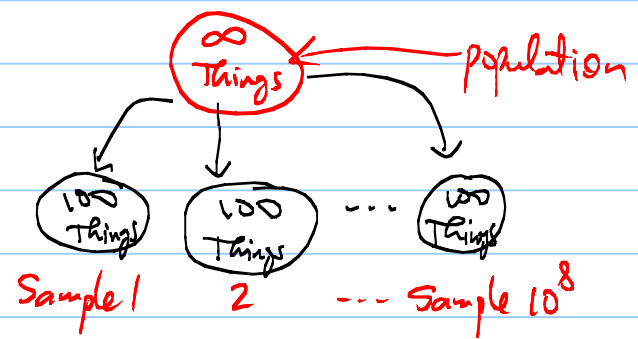
Don't confuse the various proportions: $p(X=x)$ ← Important!
 π ← prop. of 1's in each sample of size n

Irrelevant! ← It does not show-up in Binomial. It will later (Ch 7.8).

Example 1.23 (p.56)



Assume the lots are identical, i.e. the company manufacturing the 5000 things is extremely consistent. Then, the picture looks like this:



Q What proportion of these 10^8 lots will have $X=0, 1, \dots, 100$?

$G = \text{Good}$ $\text{Sample} = \{G, G, \dots, G\}$
 $B = \text{Bad}$ $\text{Sample} = \{B, B, \dots, B\}$

Suppose we know the prop. of Bads, period, in the pop. = 0.5%

Then $P(X=x) = \binom{100}{x} \pi^x (1-\pi)^{100-x} = .005 = \pi$

prop. of lots with $X=0$: $\binom{100}{0} \pi^0 (1-\pi)^{100} = .6058$

$= 1$: $\binom{100}{1} \pi^1 (1-\pi)^{100-1} = .3044$

$= 2$: $\dots = .0757$

$= 3$: Etc. $\dots = .0124$

Lab.

Important Interpretation

In the long-run we expect { ~60% of the lots to be all good.
 ~30% " " " to have 1 bad out of 100.
 ~7% " " " 2 bads " " "
 (i.e. 7% of the lots to be 2% defective)

hw-lect6-1 By R

Consider one of the two continuous variables, and one of the two discrete variables, in hw-lect 1. 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 cont. var. all on the same plot, side-by-side. Interpret!

hw-lect6-2 By R

Today, in class, a student (who we shall call Chris) asked about the effect of sample size on boxplots. I said that boxplots are generally not affected by sample size. Let's see that. Write code to

- a) take a sample of size 20 from a normal with $\mu=0$, $\sigma=10$.
- b) " " " 30 "
- c) " " " 40 "
- d) " " " 50 "
- e) " " " 100 "
- f) Make a comparative boxplot of the 5 samples a-e.

hw-lect6-3

- a) Use the binomial mass function to show that the prob. of getting "at least 1 head out of n tosses" is $1 - (1 - \pi)^n$, where π is the prob. of getting a head on a single toss.

Show work!

- b) What is the numerical value of that prob. as $n \rightarrow \infty$?
Think about the answer you get; it's interesting and counterintuitive.

This document was created with Win2PDF available at <http://www.win2pdf.com>.
The unregistered version of Win2PDF is for evaluation or non-commercial use only.
This page will not be added after purchasing Win2PDF.