

1/20/15 263 Lecture 2

- Quiz 1 on Thurs
- Website is updated w/ Excel Hw 1
- Excel Hw 1 is DUE Thurs in class or in office (M 702) Before 5:30.
- Webassign for Ch1 DUE Fri.

Last time:

1.1 - 1.3 Data Sets, Histograms, Summary Statistics.

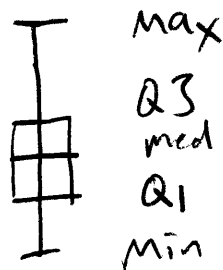
Data set: Cases \rightarrow a particular instance
Variables \rightarrow Quantities of ea. case
eg "name"
"weight"

Distribution of a Variable:

The values taken & how often they are taken.

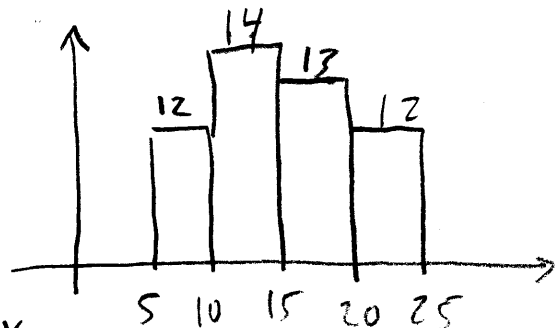
5-#-Summary:
of a distribution

1.5 IQR rule for outliers



Histogram:

Choose bins w/
equal width, &
count # of data
points with variable X
in each bin.



Mean & Std Deviation:

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N x_i$$

"ADD values & divide by number of data PTS"

Variance

$$S^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X})^2$$

Std Deviation

$$S = \sqrt{S^2}$$

"Averaged" ^{squared} DISTANCE from data PTS TO Mean"

* Why squared? generalization of Euclidean (Pythagorean) Distance, want to eliminate \pm sign.

* Why divide by $n-1$ instead of N ? Knowing the mean reduces "degrees of freedom" by 1. It's more accurate! (see Wikipedia "Bessel's correction")

* Why S instead of S^2 ? S will have the same UNITS as the variable X .

Robustness of \bar{X} & S :

Because each data PT is given "equal weight", Mean & Var/std dev are sensitive to outliers. $Q1, Q2, \dots$ are NOT!

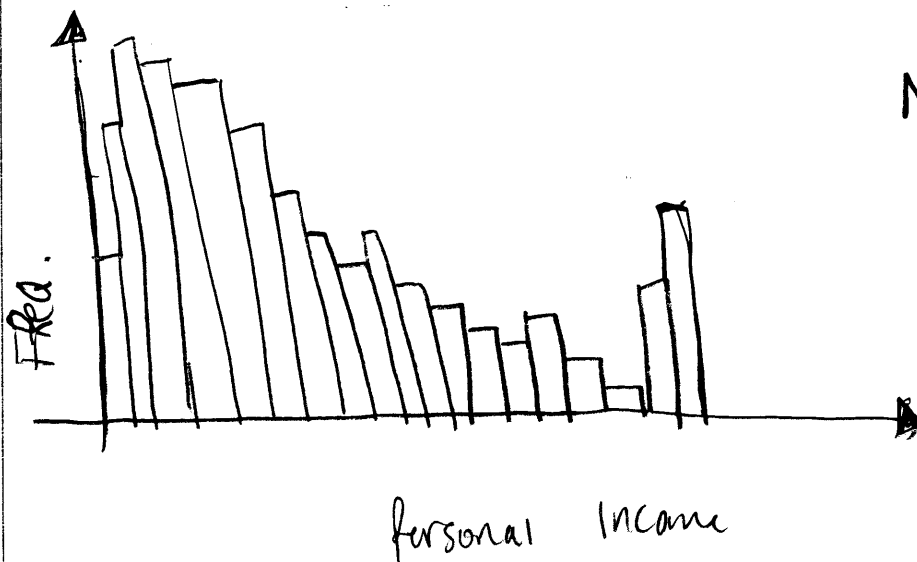
1.4: Normal Distribution, z-values

Note: Quantitative Variables can either take a discrete # of values, e.g. class size, SAT score, # of bacteria - Or - a continuous # of values i.e. exact length, sea level, temperature.

Occasionally, we will use continuous vars as a model for discrete ones. Why? Calculus. (Note: Calc is not req'd for 263!)

Again: a distribution shows the values that a variable takes on the horizontal axis, & the frequency taken by those values on the vertical axis.

EX Income Distribution, from US census.

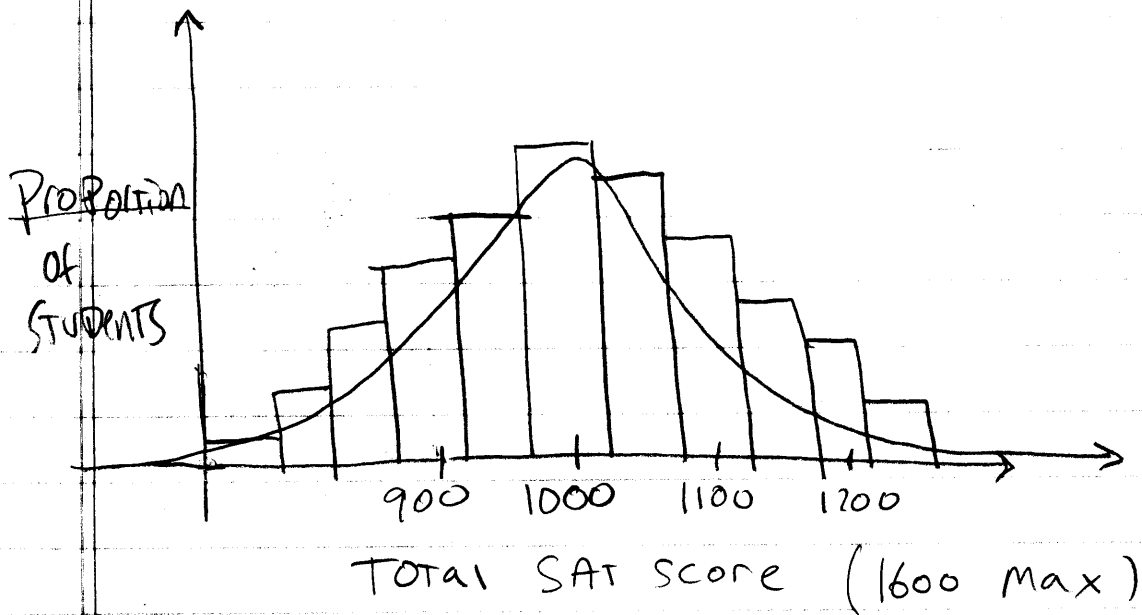


Note the shape - skewed "right" (long right tail)

Units:

- * Horizontal axis is Income
- * Vertical axis is # of people - or - Proportion of Population.

Density Curves:



Mean: $\bar{X} \approx 1000$

Std Dev: $S \approx 100$

NOTICE Density curve is symmetric, unimodal (one peak), has a mean of 1000, STD Dev ≈ 100 . "Bell-shaped"

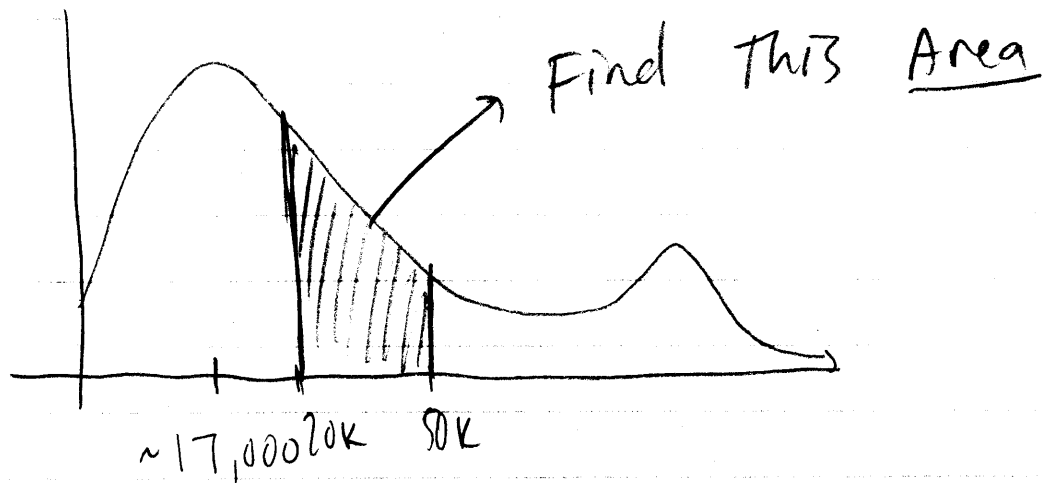
- getting a high score seems as common as getting a low score

- many scores are bunched around a single score

- a large portion of scores are concentrated w/in 2-3 std devs from \bar{X}

How will we use distributions?

- Want to answer questions like:
"What proportion of the population earns an income between \$20K & \$50K?"



Note:

For a discrete variable, we can ask questions like "How many classrooms have 32 students?"

For a continuous variable, we can only ask for an interval of values
"How many days have a temp between 57°F & 72°F ?"

For calculus-minded people: We are computing integrals of distributions.

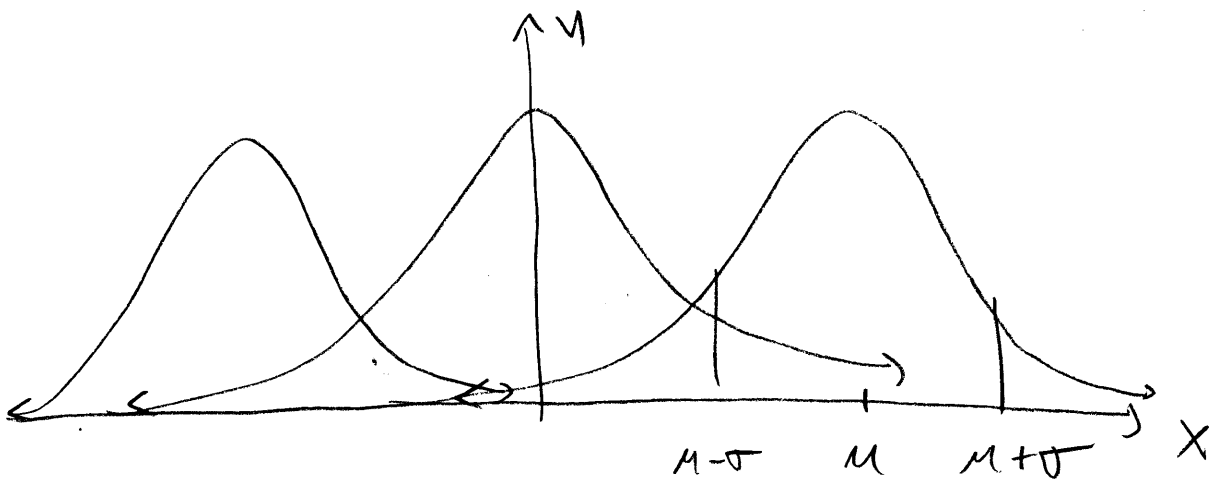
Normal distributions

A "normal" or "Gaussian" distribution is the bell-shaped curve w/ graph

$$N(\mu, \sigma) \quad y = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Two Parameters describe it: μ a mean
 σ ~ STD DEV.

Because of the central limit theorem, Normal distributions are very common. We will return to this in chapter 5.



Properties

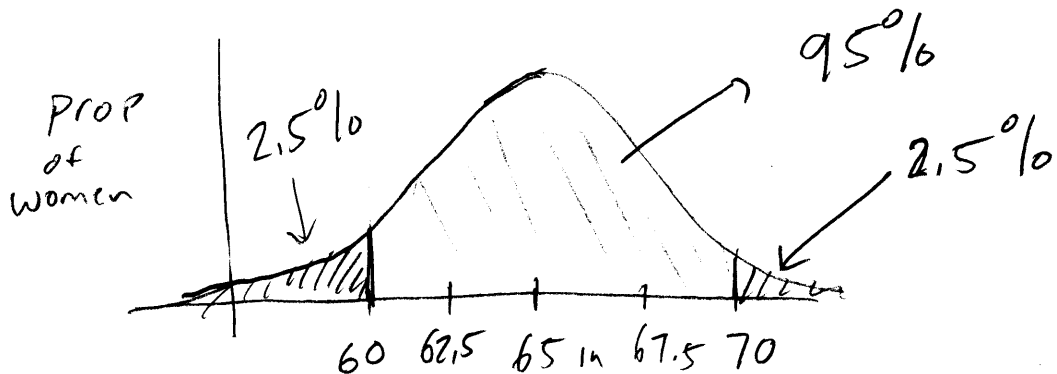
- Total area is 1 ("100%")
- 68-95-99.7 rule:
 - 68% of data is within 1 std dev of mean
 - 90% " " " " 2 " "
 - 99.7% " " " " 3 " "

Normal Distribution Examples

Suppose

Women's heights are normally distrib. w/ mean 65 inches & std dev 2.5 in.

- a) What Prop. of Women are less than 60 in?



- b) What Prop are less than 70 in?
 $95\% + 2.5\% = 97.5\%$

- c) What Prop. are more than 72 in?

Z-Scores

We "change units" so we can always use $N(0,1)$:

$$Z = \frac{X - \mu}{\sigma}$$

Units of Z are "standard deviations away from mean".

IE If $Z = 1.5$, X is 1.5 σ 's to the right of μ . If $Z = -1.6$, X is 1.6 σ 's to the left of μ .