

Histograms

10/2/12

~ defined: special type of bar graph in which the bars touch

~ width of bar has meaning (age, range, distance, time)

~ data grouped into classes

~ useful with large amounts of data

~ shape tells you info

~ Deviations hard to point out

~ exact data is hidden

~ How to choose classes

~ would info be hidden/lost?

~ will anything stand out?

~ 5-10 bars ... Mr. N likes 7

~ Width

~ $W = \frac{\text{large} - \text{small}}{\# \text{ of bars}} = \text{answer (round up!)}$

~ Basically columns to organize raw data before graphically displaying them: Frequency Table

~ Classes/groups, tallies, frequencies, midpoints, relative frequency

~ first class begins w/ smallest, next is $\text{small} + W$

~ Types of Histograms

~ Normal (symmetric)

~ Rectangular

~ Bimodal

~ Skewed left/right (tailing off side is in name)

~ Analysis (for histograms, dot plots, stem and leaf, box and whisker, etc)

~ center

~ shape

~ spread

~ outliers

~ Ex: high = 44, low = 14 $W = \frac{44-14}{10} = 3$

Classes	Marks	Tally	FREQ (%)	REL. FREQ (%)
---------	-------	-------	----------	---------------

14-16				
+3	17-19			
:				

41-44 ← include last one, to make 10 classes

~ Histograms

- ~ Bars are same width and always touch
- ~ width represents a quantitative value
- ~ height indicates frequency

For stats class:

~ don't use boundary, use limit

~ classes

~ Number: 5-15

~ width: $W = \frac{\text{Largest} - \text{smallest}}{\# \text{ of classes}}$ round up!

~ include tallies on frequency table, no boundary, no MP

~ Midpoint

~ center of class or "class mark"

~ $MP = \frac{\text{lower class limit} + \text{upper class limit}}{2}$

~ class boundary

~ Halfway point between upper limit of one class and lower limit of next class

~ used as endpoints for bars on histogram (in some conventions...)

~ Frequency Tables

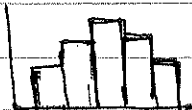
~ Determine class width \rightarrow create classes \rightarrow tally data \rightarrow obtain class frequency (total the tallies) \rightarrow compute midpoint \rightarrow determine class boundaries

~ Relative Frequency table: do above plus calculate relative frequency

$$RF = \frac{f}{n} = \frac{\text{class frequency}}{\text{total frequency}}$$

~ Distributions

~ symmetrical



~ uniform



~ skewed left



~ right



~ bimodal



Graph Analysis

10/4/12

~ Circle, line, bar, and pictographs

- ~ trend
- ~ unusual?
- ~ highs and lows
- ~ why?

~ Dot plot, Histogram, stemplot, box plot

- ~ center
- ~ shape
- ~ spread
- ~ outliers

Stem and leaf plots

- ~ easily constructed
- ~ shows original, specific data
- ~ shows shape and distribution
- ~ allows analysts to view...

~ small quantities of data

~ How to:

- 1) Find smallest and largest value
- 2) Decide on stems and write them vertically w/ a line to right
- 3) separate values into stem and leaves. Put leaves on right
- 4) on a new plot, arrange leaves from smallest to largest
- 5) Have a key (ex: 2|3 = 23), title, make sure #'s line up

~ Analysis

- ~ center: median, be sure to indicate what kind; ex: 22 (median)
- ~ shape: turn graph to side \rightarrow ome. it's a histogram
- ~ spread: where is the bulk; number and word
(ex: 12-24 (range of 16); moderate)
- ~ outliers: what is far away from the bulk

~ Example:

- 1) Domestic (mpg): 23, 24, 15
Foreign (mpg): 33, 35, 27

if you have a lot of data:

Domestic		Foreign		Dom.		FOREIGN
5	1			5	1	1 → 10-14
4	3	2	7	4	3	• → 15-19
3	3	5	3/3=33	•	7	
3/2=23	4			3	3	
				•	5	
CSSO		CSSO		4		
				•		
				CSSO		CSSO

2) From 2.3 #8

1	5
2	
3	
4	9
5	4
6	
7	
8	5
9	0 5
10	0 2 2 6
11	
12	3 6
13	6 0 9
14	4 9
15	0 4 9
16	6 3
17	5
18	5
19	
20	
21	
22	
23	5

3) how off are you?

	women	man
	36	28
	4	8
	12	20
	3	40
	2	20
	13	12
	20	15
	18	5
		10

23/5 = 23.5

Section 2.3

~ Exploratory Data Analysis (EDA)

~ useful for detecting patterns and extreme data values

~ Stem and Leaf Display

~ one EDA technique

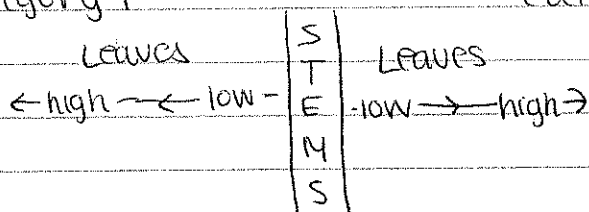
~ organizes and groups data, but allows us to see original data too

~ How to: (see class notes)

~ Back-to-Back Stem and Leaf Plots

category 1

category 2



Against All Odds video 2 Notes

10/8/12

~ Histogram (Lightning Example)

~ Distribution: overall pattern (first flash b/w 11-12 noon, symmetric)

~ Outliers: data that is away from the majority

~ Good for large quantities of data

~ Always look for the big picture first

~ Shape, Center, Distribution/Deviation

~ Symmetry, skewed left/right, bimodal

~ Center: where is the tallest bar, median

~ Outlier: you want to be able to explain why? ~~the~~

~ spread

~ If spread is too much, graph is less useful → split it up

~ Stemplot

~ Center: median

~ Outlier: why?

~ Keep exact data, good for small quantities of data

~ CSSO, then conjecture why

Semester one Project (due Oct 29)

~ Procedure

~ script

~ so that each partner says same thing

~ have the person step away

~ Visual, Auditory, Motion (for assurance about bag and incentive)

~ BIG BAG. with surveys in it always! and incentive!

~ sample size

~ 100 or more

~ MUST GO OFF CAMPUS AS WELL (at Medea, check-in day before, approval)

~ subgroups must be 25 or greater

~ Survey

~ Instructions: remind of bag and incentive

~ 12-15 questions, one page

~ Read tutorials

~ 5 peer reviews

~ 13 kids to convenience sample as practice/test run

~ Run it by Mr. Micek

Measures of central tendencies

10/10/12

~ The Averages

~ Mean

~ Median

~ Mode

~ Trimmed

~ Arithmetic Mean

~ usual meaning behind "average"

~ $\text{mean} = \frac{\text{sum}}{\# \text{ of items}}$

of items

~ Population = ΣM

~ Sample = \bar{X}

~ Outliers can have large impact. Pulls mean toward outliers

~ Median

~ middle value when ordered small to large

~ position, not numeric values \rightarrow resistant to outliers

↓

~ ex: 13 15 17) (18 19 20 median = 17.5
 13 15 17) (18 320 450 median = 17.5 ← should it be?

~ mode

- ~ most often occurring value
- ~ overlooked, but often more appropriate
 - ~ ex: average hat size
- ~ can be bimodal or trimodal or multimodal; or no mode
- ~ outliers have no effect

~ Trimmed Means

- ~ mean that resists extremes; eliminates pull of extremely low or high values
- ~ n% trimmed mean → top of n% off each end (Round up!) → find mean as normal

Bird Sneak

- ~ mean, median, mode, 5% Trimmed mean (show work)
- ~ Pattern?
- ~ Use any wild animal (use same one every time, 10 keys)
- ~ Prediction, Procedure, Data, Graphs (w/ CSSO), conclude

Measures of Variation

10/12/12

- ~ Population Parameters: μ, σ^2, σ
- ~ Sample Stats: \bar{x}, s^2, s
- ~ Standard deviation

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$

~ Variance

- ~ standard deviation squared

~ examples

~ TEMP. in last 10 years

-8 x

3 x

25 x

17 x

-2 x

10 x

-12 x

21 x

4

6 s

$$\bar{X} = 6.4$$

$$\text{median} = 5$$

$$\text{Range} = 25 - (-12) = 37$$

$$s^2 = 146.49$$

$$s = 12.10$$

~ Ruby weight

19.8

43.8

36.1

52.4

63.1

20.7

46.3

$$\text{range} = 43.3$$

$$\mu = 40.31$$

$$\sigma = 14.82$$

Sections 3.1 and 3.2

~ Mode: most common value

~ Median: odd (middle number once small \rightarrow large); even (average of the two middle numbers once small \rightarrow large)

~ Mean: sum / # of entries

~ sample mean: \bar{x}

~ Population mean: μ

~ Trimmed Mean: order small \rightarrow large, take x% of front and back, take mean of the rest

~ Measures of variation

~ Range: large - small

~ Standard Deviation: how data differs from the mean; small sd = small spread

~ sample standard deviation

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

* $\sum (x - \bar{x})^2 = SS_x = \text{sum of squares}$

~ Population standard deviation

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$

* uses "N" not "n-1" b/c samples often don't have extremes. "n-1" \rightarrow s bigger \rightarrow accounts for this

~ variance: standard deviation squared

~ sample variance

$$s^2 = \frac{\sum (x - \bar{x})^2}{n-1}$$

~ Population variance

$$\sigma^2 = \frac{\sum (x - \mu)^2}{N}$$

~ Coefficient of Variation: expresses standard deviation as a percent of sample or population mean (good for comparison b/c no units)

$$\text{Sample: } CV = \frac{s}{\bar{x}} \times 100$$

$$\text{Population: } CV = \frac{\sigma}{\mu} \times 100$$

~ Chebyshev's Theorem: for any set of data and for any constant $K > 1$, the proportion of data that must lie within K standard deviations on either side of the mean is at least $1 - \frac{1}{K^2}$

~ for any set of data:

~ at least 75% of data fall into interval $\mu - 2\sigma$ to $\mu + 2\sigma$

~ at least 88.9% of data fall into interval $\mu - 3\sigma$ to $\mu + 3\sigma$

~ at least 93.8% of data fall into interval $\mu - 4\sigma$ to $\mu + 4\sigma$

~ Spread (SSO)

~ use coefficient of variation

~ CV

word

↓ 10-13%

Tight

15-30%

Low

33-67%

Medium

70-95%

High

97%↑

Wild

} within each range
use "moderately"
or "very"

~ ex) spread: 10.28% CV, Tight

~ Grouped data (ex: frequency table)

$$\bar{x} = \frac{\sum xf}{n}$$

* x = midpoint of the class

$$n = \sum f$$

$$s = \sqrt{\frac{\sum (x - \bar{x})^2 f}{n - 1}}$$

f = # of entries in the class

$$s = \sqrt{\frac{SS_x}{n - 1}}$$

$$SS_x = \frac{\sum x^2 f - (\sum xf)^2}{n}$$

Box and whisker plots

10/18/12

~ Percentile

~ like mile markers; 10th percentile = above/equal to 10% of people

~ There is no 100th percentile, highest is 99th percentile

~ ex: P. 112 #3

Not fair. You would fire 70% of staff. If all scores were very high, 82% could be in the bottom

~ Quartiles

~ special percentiles that splits data into fourths

~ 1st quartile = LQ = Q_1 = 25th percentile

~ 2nd quartile = median = Q_2 = 50th percentile

~ 3rd quartile = UQ = Q_3 = 75th percentile

~ 4th quartile = ... = 100th percentile

Section 3.3

~ Percentile

~ for whole numbers P (where $1 \leq P \leq 99$), the P^{th} percentile is a value such that $P\%$ of the data fall at or below it and $(100 - P\%)$ of the data fall at or above it.

~ QUANTILES

~	25%	25%	25%	25%
	Q_1	Q_2	Q_3	
	LQ	Median	UQ	
	25%ile	50%ile	75%ile	

~ TO COMPUTE:

- 1) order small \rightarrow large
- 2) Find Median - this is Q_2
- 3) Q_1 is median of data below Q_2
- 4) Q_3 is median of data above Q_2

~ Inter Quartile Range = $Q_3 - Q_1$

~ 5 Number summary

~ min, Q_1 , Q_2 , Q_3 , max

~ BOX and whisker

~ find the 5 Number summary $\hat{=}$ IQR

~ Draw a box around IQR

~ Draw a line through median

~ Draw whiskers out to min & max

~ TITLE & AXES LABELED

~ CSSO!!

