

- * Studying data: formalizes process
- * Statistics - describes any bit of information
- * 1986

Challenger - wings did not seal bc cause of cold

- plot ranges are very off - only showed bad points
- resilience of rubber was affected by temperature

Population: all potential observations from a distribution of interest

enumerative study:

- tangible population
- sampling frame: list of population exists or possible
- finite

Analytic study:

- ongoing process
- conceptual/hypothetical
- ~ infinite ex: rainfall over time
ex: burning pop. - an 's gallon milk container

Sample:

measurements ex: height, etc
make up dataset

Observational Study

- scientists cannot control variables of interest
- step after controlled studies
- is sample representative of pop?
- census or complete enumeration all in computing - every person is measured
- + random sample is best
"simple random sample
draw names from hat"
- Sample variation:
sample tends to be diff from population
- samples of convenience:
- can be biased
leads to error

Variation:
- people choose to part or not

* historical control

Variable: characteristic we are interested in

Data

Categorical - not/partially additive

numerical - additive

univariate - one info bit

bivariate - 2 info bits

tri/quadrivariate (etc)

ex: $x = (h, c)$

descriptive stats:

describes + simplifies data sets

inferential statistics:

attempts to say something about data

ROT PLOTS

univariate data

question typical datapoint with like...

Histograms

- bar chart for numbers
- shape describes distribution

Create

- 1) find range (max - min)
- 2) Break into intervals
 - ~ sample size
 - class = range / # classes
- 3) freq. chart
- 4) column chart

ex: range = 44.6

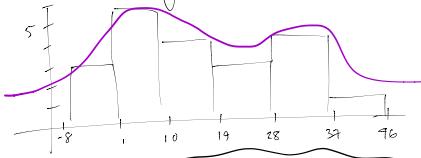
5 classes

class size = $\frac{44.6}{5} = \sim 9$

class	freq	relative freq.
[-8, 1)	3	0.15
[1, 10)	5	0.25
[10, 19)	4	0.2
[19, 28)	3	0.15
[28, 37)	4	0.2
[37, 46)	1	0.05
(total data)		= 1
$\frac{44.6}{10} = 4.46$		

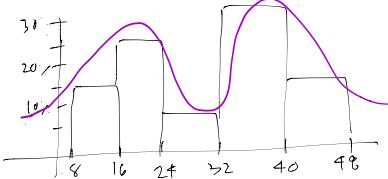
how do you decide class size?

histogram



ex 2 6 classes
 $\frac{44.6}{6} \sim 8$

class	freq	R freq
(-8, 0)	3	15
[0, 8)	5	25
[8, 16)	1	5
[16, 24)	6	3
[24, 32)	2	1
[32, 40)	3	15



Shapes

- dist. left skewed neg.
- symm
- right skewed positive
- bimodal
- multimodal

* common platines are more trustworthy

* shapes related to real life events

Density Histograms

use for pre-binned data

Measures of location

- summary of data set
- sample statistic - derived from Sample
- sample mean - most common for sample center
- μ - mean from all population, not sample
- mean is greatly effected by outliers

outlier very different from rest, generally wrong data

robust resistant to outliers

(mean = \bar{x}) $x_{(i)} = \text{difference}$
(median = \tilde{x}) $x_{(i)} = \text{smallest, etc}$
median - order stat, more robust

If n is odd, $\tilde{x} = x_{(\frac{n+1}{2})}$

If n is even, $\tilde{x} = \frac{x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}}{2}$

$\tilde{\mu} = \text{median of population}$

categorical data - measured in proportions (from population)

$$\hat{p} = \frac{x}{n}$$

+ skew = $\bar{x} > \tilde{x}$

- skew = $\tilde{x} > \bar{x}$

Symm = $\bar{x} = \tilde{x}$

Variability/spread

range \rightarrow susceptible to outliers

deviations from mean

how far from mean for each datapoint

- + - dev. cancel each other out, = 0

sample variance: measures spread of data

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} S_{xx}$$

sample standard deviation:

$\sqrt{}$ of sample variance

$$s = \sqrt{\frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right]}$$

e.g. 19 data points

$$\frac{1}{19} \left[\sum_{i=1}^{20} x_i^2 - \frac{1}{20} \left(\sum_{i=1}^{20} x_i \right)^2 \right]$$

$$S^2 = 186.49$$

$$\sqrt{186.49} = 13.65 = s$$

NOT robust

$\sim 2/3$ data falls w/ 1 SD

$\sim 95\%$ in 2 devs

\sim all data in 3 std
 $(\bar{x} - 3s, \bar{x} + 3s)$

true for normal distribt.

$$\begin{aligned} \text{ex } (\bar{x}-s, \bar{x}+s) &= (1.71, 29.02) \\ (\bar{x}-2s, \bar{x}+2s) &= (-11.95, 42.68) \end{aligned}$$

Ipad died @ January 25, 2016

Question 8:

Average = sum of X_i 's/10

$10(63,000) = \text{sum of } X_i$'s

$\frac{1}{10}$ th of old number subtracted
add another $\frac{1}{10}$ th of new number, = new total

$n-1$ in variance - corrects for bias

n in population variance

usually don't know population mean, only have sample

Finding Quartiles

Q1: median of bottom half:

Q2: median

Q3: median of top half:

-robust measure of spread

-interquartile range or fourth spread: Q3-Q1

Boxplot

Plot lower quartile and upper quartile, put box around it and a line where median is

"whiskers" where greatest and least last non-outliers exist

outliers: inner quartile range * 1.5

if no outliers exist, the whiskers sit on top and bottom values

Inner quartile range = Q3-Q1

small outlier cutoff: $Q1 - 1.5 * (\text{IQR})$

large outlier cutoff: $Q3 + 1.5 * (\text{IQR})$

DIFFERENT THAN WHISKER POINTS

Right skewed



$\therefore \text{mean} > \text{median}$
(usually)

left



$\therefore \text{mean} < \text{median}$

Not clear if skewed or not if whisker is out further, for example:



$$\text{IQR} = Q3 - Q1$$

$$\text{Ex: } 72.6 - 21 = 51.6$$

outliers, small

$$21 - 1.5 * 51.6 \\ - 56.4$$

$$72.6 + 1.5 * 51.6 \\ 150$$

Histogram vs. Box Plot?

Histogram shares more information, but difficult to compare

Box Plots are easier to compare

Probability:

Dealing with chance, randomness, and uncertainty

Experiment:

Cannot be determined in advance what outcome will be

Sample Space:

Set of all possible outcomes

Randomly pull card?

Cards:

$S = \{\text{ace of spades}, 2 \text{ of spades}, \text{ace of hearts, etc...}\}$

Event: draw a hearts (subset, A-K hearts)

Component breaks on computer?

$S = (0, \infty)$ or $[0, \infty)$ or $[0, 1 \text{ million}]$

Subset = $[0, 20]$

Event: "interesting" subset of sample space

Empty set – subset of every set

Sample range

highest - lowest

Sample standard deviation:

Sum of squares of differences
divided by $(n-1)$ then $\sqrt{\dots}$

of finer answer !!