

Housekeeping.

- Exam 1 is Friday (!)
- Optional homework posted - sol'ns to be posted tonight
- No written homework for Friday (just study...)
- Week 2 reading summary + discussion due Monday 11:59 p.m.
- ~~✱~~ No class on Monday.

# Study Guide.

2/21

- Population / sample
- Parameter / Statistic
- Variables / data
- Descriptive / inferential statistics

- Qualitative data

vs.



- Dot plots
- Freq. tables
- Histograms
- Stem & leaf diags.
- What makes a good infographic?

- Nightingale — military deaths in Crimea
- Minard — Napoleon march.

• Sampling

• Simple random

• Stratified

• Cluster

• Opportunity

• Systematic

technical problems - <sup>e.g.:</sup> how to select proportionally the sample candidates from each stratum

• Plotting data in pie charts

- Sorted pie charts

• ————— bar charts

vs. unsorted

- Bar charts vs.

Pareto charts.

• Using random # generator.

• Infographic: Premier League Results (Financial Times)

• Ethical issues w/sample selec'n.

• When to use (or not) certain sampling techniques

- Costly, controversial, impossible situations to create sampling frame

- Cluster sampling is cheaper than stratified sampling

- Problems w/opportunity (a.k.a. convenience) sampling.

• Infographic: John Snow Cholera map ? Voronoi diagram.

- 13  
4
- Sampling error
  - Non sampling error
  - Sample bias
- } real-world examples.

- Sampling  $\left\{ \begin{array}{l} \text{with} \\ \text{vs.} \\ \text{without} \end{array} \right.$  replacement.

- Mean, median, mode : measures of center of data

- Outliers : the mean is much more sensitive to outliers than the median is.

- Distinction btwn. measures of ctr.  $\hat{=}$  meas. of spread.

- Creating a freq. table or histogram — determining class bdris.

- Infographics : scatter plot of temps. (NYT)

time series of global temps. (xkcd.com)

• Descriptive statistics

Sec's (from O.S. book):

2.1

2.2

2.3

2.4

2.5

2.6

2.7

• Box plots (box & whisker plots)

• Quartiles & IQR

• Std. deviation & variance } opt HW assignment

• z-scores

• Infor. (not totally relevant): corn / milho.

S.D.  $\hat{=}$  z-scores.

Measures of center - mean, median, mode  
- tells the "typical" value of a pt. in a data set

Measures of spread - IQR, s.d., variance  
- tells how spread out the data is about its central value.

IQR measures how spread out the data is abt. the median;

S.D.  $\hat{=}$  variance  $\longrightarrow$  mean.

Naïve idea/ Find the average distance btwn. <sup>data</sup> points  $\hat{=}$  the mean:

First try:

- Compute the mean, call it  $\bar{x}$ .
- For each data pt. in the set, compute the distance btwn.  $\bar{x}$   $\hat{=}$  that point.  
i.e., if  $x_i$  is a pt. in the set, compute  $\bar{x} - x_i$ .
- ~~Calculate~~ Find the mean of those distances.

Try this algorithm out for the set  $\{3, 6, 4, 3\}$

•  $\bar{x} = \frac{3+6+4+3}{4} = \frac{16}{4} = 4.$

• Distances =  $\{4-3, 4-6, 4-4, 4-3\} = \{1, -2, 0, 1\}$

• Spread =  $\frac{1 + (-2) + 0 + 1}{4} = \frac{0}{4} = 0.$

Slight modification:

• Compute  $\bar{x}$

• For ea. data pt., compute  $(\bar{x} - x_i)^2$ .

For technical reasons, when computing this mean, divide by  $n-1$ , not  $n$ , where  $n$  is the # of data pts. in the set.

• Find the mean of those, and call it the variance, denote it as  $s^2$ .

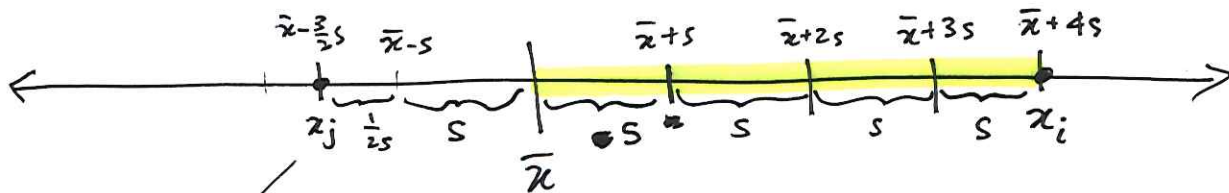
• Take the sqrt of the variance, call it the s.d.,  $s$ .

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

• s.d. is a property of a data set.

• z-score  $\rightarrow$  point in a set.

z-score of a pt. tells you how many s.d.'s ~~are~~ above the mean that pt. is.



z-score of  $x_j$  is  $-\frac{3}{2}$

z-score of  $x_i$  is 4

Here the z-score of  $x_i$  in a set w/ mean  $\bar{x}$   
and s.d.  $s$  is:

$$z := \frac{x_i - \bar{x}}{s} .$$