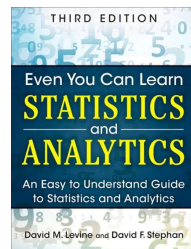


IMGD 2905

Descriptive Statistics

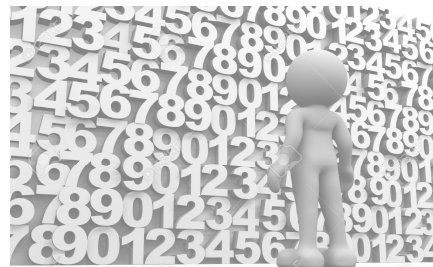
Chapter 3



1

Summarizing Data

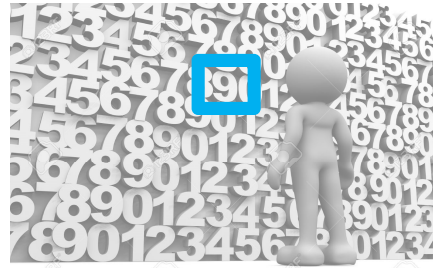
- With lots of playtesting, there is a lot of data
 - This is a good thing!
- But raw data is just a pile of numbers
 - Rarely of interest
 - Or even sensible
- **Q:** How to summarize all this information?



2

Summarizing Data

- With lots of playtesting, there is a lot of data
 - This is a good thing!
- But raw data is just a pile of numbers
 - Rarely of interest
 - Or even sensible
- **Q:** How to summarize all this information?



Measures of **central tendency**

Examples?

3

Measure of Central Tendency: **Mean**

The sum of the measurements

divided by the number of measurements

$$(6 + 4 + 5 + 4 + 8 + 3) / 6 = 5.$$

gives you the mean.

<http://www.cdn.sciencebuddies.org/Files/463/9/MeanEquation.jpg>

- Also called the “**arithmetic mean**” or “**average**”
- In Excel, =AVERAGE(range)
=AVERAGEIF() – averages if numbers meet certain condition



4

Measure of Central Tendency: Median

- Sort values low to high and take middle value



10 11 13 15 16 23 26

middle number
<https://www.mathsisfun.com/definitions/images/median.gif>

13 22 26 38 36 42 49 50 77 81 98 110

Median = 45.5
<http://www.nedarc.org/statisticalhelp/basicStatistics/measuresOfCenter/images/median.gif>

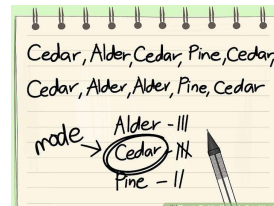
- In Excel, =MEDIAN(range)



5

Measure of Central Tendency: Mode

- Number which occurs most frequently
 - Not too useful in many cases
- Best use for **categorical data**
- e.g., most popular Hero group in Heroes of the Storm



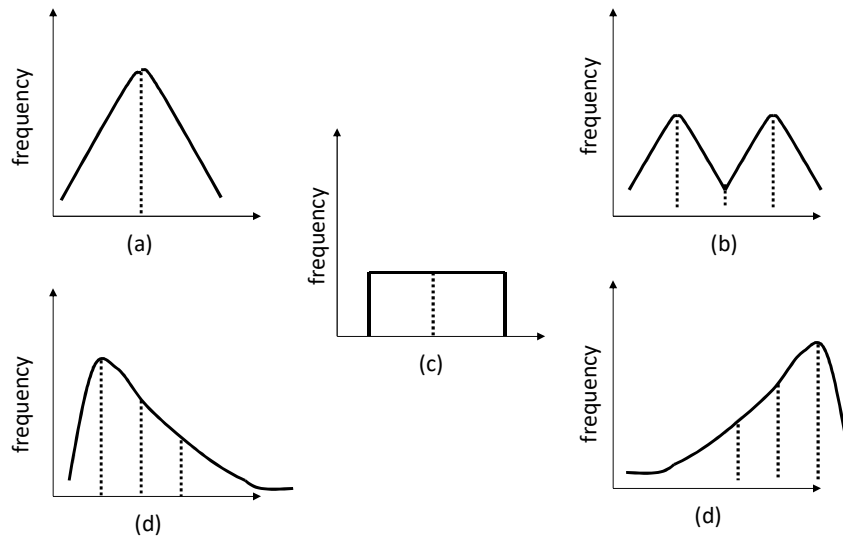
<http://pad3.whstatic.com/images/thumb/cd/Find-the-Mode-of-a-Set-of-Numbers-Step-7.jpg/aid130521-v4-728px-Find-the-Mode-of-a-Set-of-Numbers-Step-7.jpg>

- In Excel, =MODE()



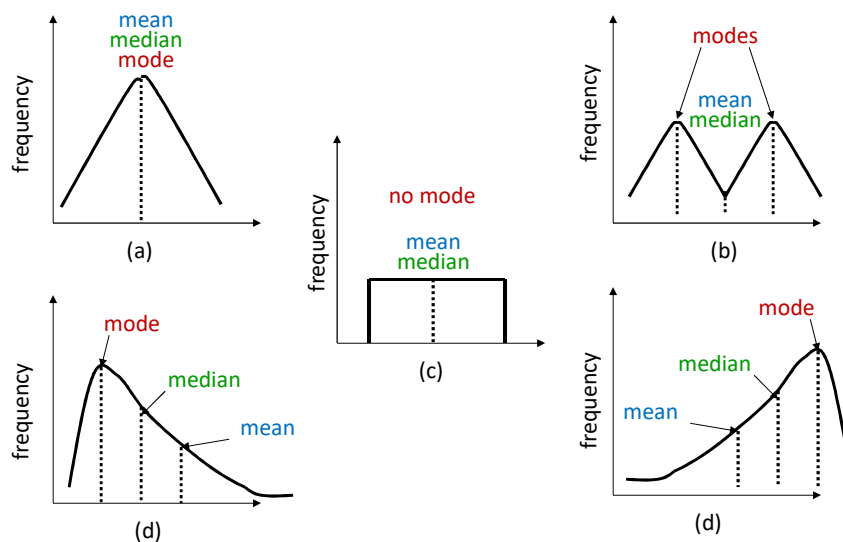
6

Depiction: Mean, Median, Mode?



7

Depiction: Mean, Median, Mode?



8

Which to Use, Mean, Median, Mode?

9

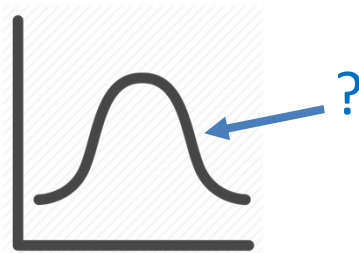
Which to Use, Mean, Median, Mode?

- **Mean** many statistical tests with sample
 - Estimator of population mean
 - Uses all data
- **Median** is useful for skewed data
 - e.g., income data (US Census) or housing prices (Zillow)
 - e.g., *Overwatch* team (6 players): 5 people level 5, 1 person level 275
 - Mean is 50 - not so useful since no one at this level
 - Median is 5 - more representative
 - Does not use all data. “Resistant” to extremes (e.g., 275)
 - But what if were exam scores? Hard to “bring up” grade
- **Mode** is useful primarily for categorical data only
 - Most played League champion, most popular maze, ...

10

Other Measures of Position

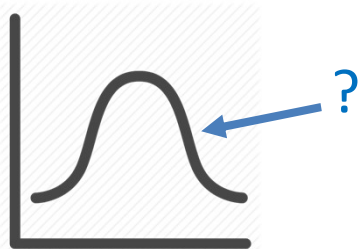
- May not always want center
 - e.g., want to know best League Champions
- What other positions may be desired?



11


Other Measures of Position

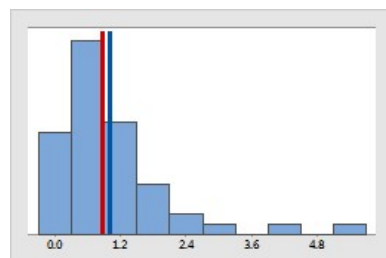
- May not always want center
 - e.g., want to know best League Champions
- **Maximum / Minimum**
 - Not discussed more
- **Trimmed Mean**
- **Quartiles**
- **Percentiles**



12

Trimmed Mean

- Take “trimming” off top and bottom (typically 5% or 10%)
 - Reduces effects of extreme values, like median
- In Excel, =TRIMMEAN(array, percent) 




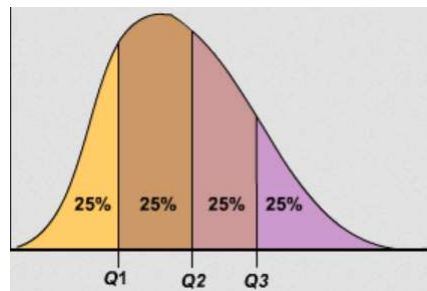
Blue – original mean
Red – trimmed mean

http://support.minitab.com/en-us/minitab/17/histogram_mean_vs_trimmed_mean.png

13

Quartiles

- Sort values
- First quartile (Q1) is 25% from bottom
- Third quartile (Q3) is 75% from bottom
- (What is second quartile?)
- In Excel, =QUARTILE(array, n) 




<https://www.hackmath.net/images/quartiles.png>

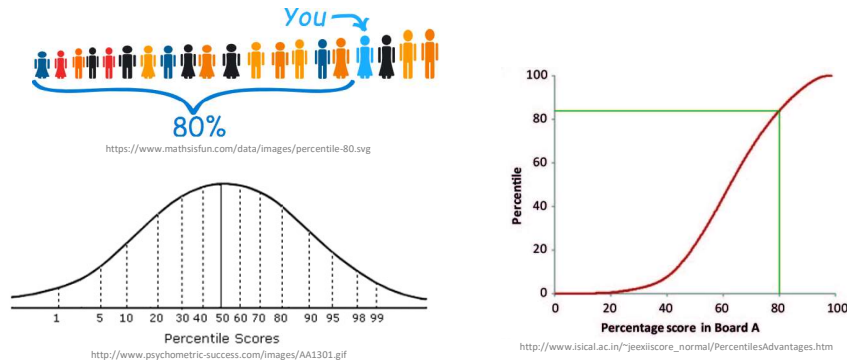
First Quartile	Second Quartile	Third Quartile	Fourth Quartile
24, 25, 26	27, 30, 32	40, 44, 50	52, 55, 57
Q ₁	Q ₂	Q ₃	
26½	36	51	

<https://mathbitsnotebook.com/Algebra1/StatisticsData/quartileboxview2.png>

14

Percentiles

- Generalization of quartiles
- N^{th} percentile is data point $n\%$ from bottom of data
- Interpolate as for first quartile 
- In Excel, =PERCENTILE(array, k) (k: 0 to 1)



15

Summarizing Data, Part 2

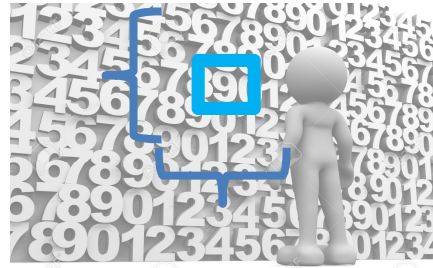
- Ok, pile of numbers can now be summarized as *one* number
 - Mean, median, mode
- But is that enough?
- Q: What other major aspect of numbers haven't we summarized?



16

Summarizing Data, Part 2

- Ok, pile of numbers can now be summarized as *one* number
 - Mean, median, mode
- But is that enough?
- Q: What other major aspect of numbers haven't we summarized?



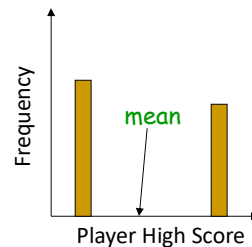
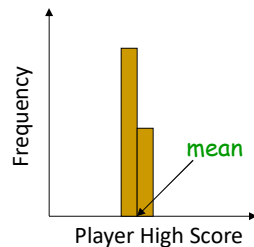
Measures of **variation**
(aka measures of *dispersion*, or
measures of *spread*)

17

Summarizing Data, Part 2

"Then there is the man who drowned crossing a stream with an average depth of six inches." – W.I.E. Gates

- Summarizing by single number rarely enough → need statement about **variation**

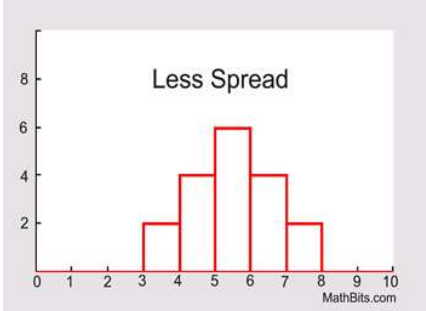


Above: does single number (**mean**) tell you enough about data?

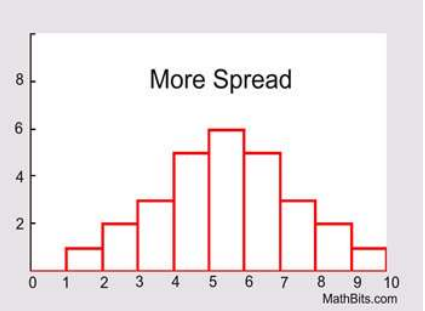
18

Variation Overview (1 of 3)

- Is data clumped or spread out?



Less Spread



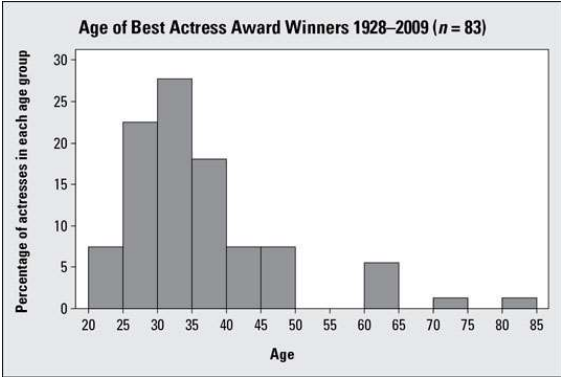
More Spread

<https://mathbitsnotebook.com/Algebra1/StatisticsData/STSpread.html>

19

Variation Overview (2 of 3)

- Is data clumped or spread out?

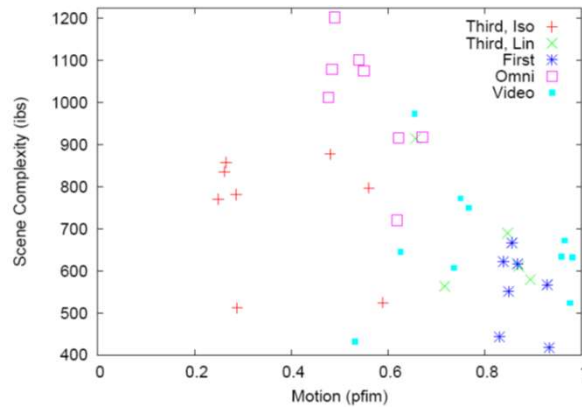


Age of Best Actress Award Winners 1928-2009 ($n = 83$)

20

Variation Overview (3 of 3)

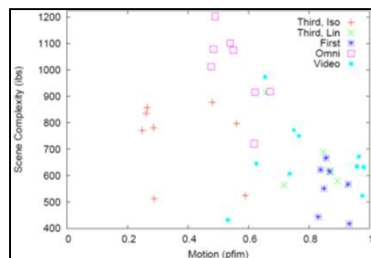
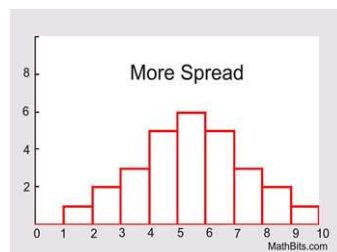
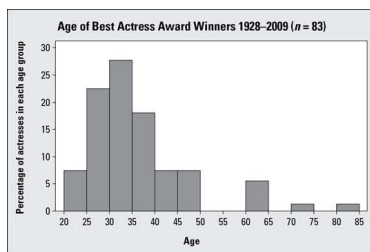
- Is data clumped or spread out?



“[Motion and Scene Complexity for Streaming Video Games](#)”

21

What are Some Measures of Variation?



22

Range

- Difference between smallest and largest value
- Somewhat obvious, but doesn't tell you much about "clumping"
 - Minimum may be zero
 - Maximum can be from outlier
 - Event not related to phenomena studied (e.g., 0 on project)
 - Maximum gets larger with # samples, so no "stable" point

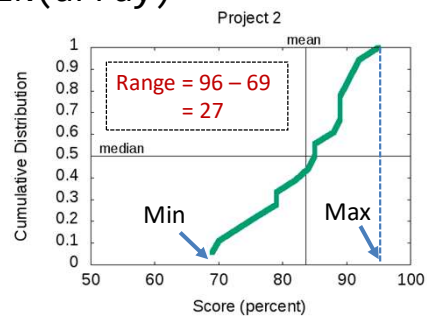


In Excel, =MAX(array) - MIN(array)

12, 25, 27, 29, 36, 38, 40, 43, 50, 54, 62

Range = 62 - 12 = 50

<http://idolosol.com/images/range-3.jpg>



23

Variance

- Compute **mean** of sample
- Compute how far each value in sample is from **mean**
 - Some can be less than **mean**, some greater
 - So square this difference (why square?)
- Divide by number of sample values – 1
 - The “-1” corrects “bias” when trying to estimate *population variance* using *sample variance*

$$\text{Sample Variance} = s^2 = \frac{\overset{\text{"sum up all"}}{\sum} (X - \overset{\text{"mean"}}{\bar{X}})^2}{n - 1}$$

24

Variance Example

- Sample kills in *League of Legends* match
 - 12, 20, 16, 18, 19
 - What is sample variance?
- First, **mean** = $85 / 5 = 17$

Kills	$X - \text{mean}$	$(X - \text{mean})^2$
12	-5	25
20	3	9
16	-1	1
18	1	1
19	2	4

$$s^2 = (25 + 9 + 1 + 1 + 4) / (5 - 1) = 40 / 4 = 10 \text{ kills squared}$$

- In Excel, =VAR(array) 

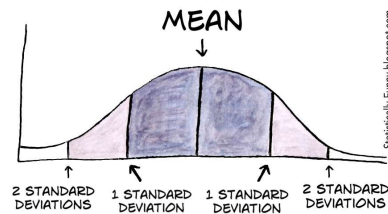
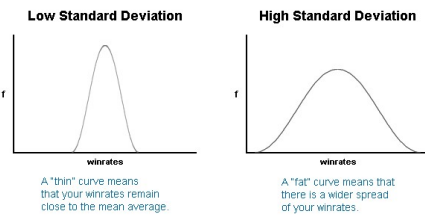
“Larger” means
“more spread”
... but units odd

25

Standard Deviation

- **Square-root of variance**
- Usually, use standard deviation instead of variance
 - Why? → Same *units* as data (e.g., “kills” in previous example)
- Can compare standard deviation to mean (*coefficient of variation*, next)
- But first:
 - **Mendenhall’s Empirical Rule**
 - **Z-score**

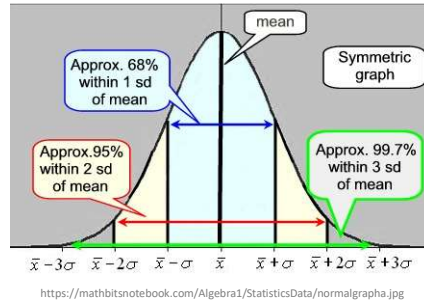
$$s = \sqrt{\frac{\sum(x_i - \bar{x})^2}{n - 1}}$$



26

Mendenhall's Empirical Rule

- About **68%** data within one standard deviation of **mean**
 - interval between **mean-s** and **mean+s** contains about 68% of data
- About **95%** within 2 standard deviations of mean
- Almost all** data within 3 standard deviations of mean



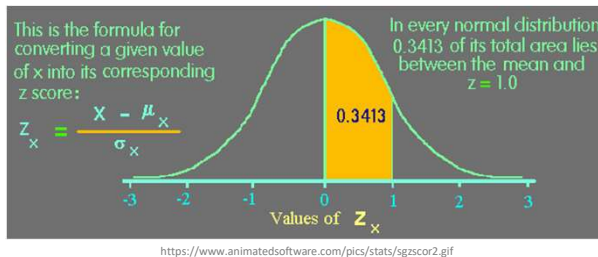
Rule assumes normal ("Bell curve") distribution

27

Z-Score

- Measure of how "far" from center (**mean**) single data point is
 - Not measure of dispersion for whole data set

$$Z = \frac{X - \bar{X}}{S}$$



Example

Mean	469
Std dev	119
X	650
Z-score for X?	
	$(650 - 469)/119$ 1.52

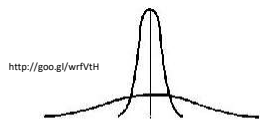
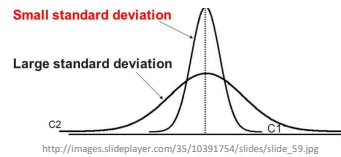
28

Coefficient of Variation (CV)

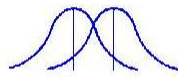
- Size of **standard deviation** relative to **mean**
 - e.g., large sd & large mean, not so spread
 - but large sd & small mean, more spread
- **Standard deviation** divided by **mean**
 - Can do this since same units!
- CV is “unit-less”, so measure of spread independent of quantity
 - E.g. seconds, clicks, spaces

Shown as percent (multiply by 100)

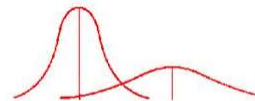
$$CV = \frac{s}{\bar{x}} \times 100$$



Same Means
Different Standard Deviations



Different Means
Same Standard Deviations

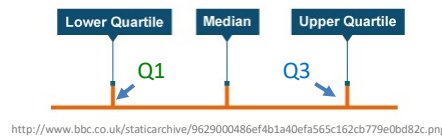


Different Means
Different Standard Deviations

29

Semi-Interquartile Range

- ½ distance between **Q3** (75th percentile) and **Q1** (25th percentile)



$$\frac{Q3 - Q1}{2}$$

- Guideline: use semi-interquartile (SIQR) for index of dispersion whenever using **median** as index of central tendency

30

Index of Variation Example

(sorted)
Lap Times

1.9
2.7
3.9
4.1
4.2
4.2
4.4
4.5
4.5
4.8
4.9
5.1
5.1
5.3
5.6
5.9

- First, sort. Then, compute:
 - Mean = 4.4
 - Min = 1.9, Max = 5.9
 - Median = $[16 / 2] = 8^{\text{th}} = 4.5$
 - Q1 = $16 / 4 = 8^{\text{th}} = 4.1$
 - Q3 = $3 * 16 / 4 = 12^{\text{th}} = 5.1$
- $SIQR = (Q3 - Q1) / 2 = 0.5$
- $Variance = 0.96$
- $Stddev = 0.98$
- $CV = stddev/mean = 0.22$
- $Range = \max - \min = 4$

31

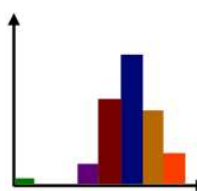
Ranking of Affect by Outliers?

Measure of Variation

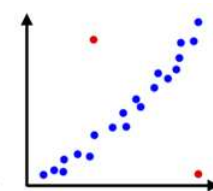
- Variance
- Range
- Standard Deviation
- Coefficient of Variation
- Semi-interquartile Range

Most to Least

?



outlier result(**green**)



outlier points(**red**)

<http://www.a-levelmathstutor.com/images/statistics/outliers-graph01.jpg>

32

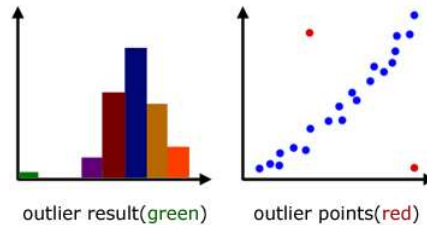
Ranking of Affect by Outliers?

Measure of Variation

- Variance
- Range
- Standard Deviation
- Coefficient of Variation
- Semi-interquartile Range

Most to Least

- Range **susceptible**
- Variance
 - Standard Deviation
 - Coefficient of Variation
- SIQR **resistant**



33

Index of Variation Summary

- Ranking of affect by outliers
 - Range **susceptible**
 - Variance
 - Standard deviation
 - Coefficient of variation
 - Semi-interquartile range **resistant**
- Note, all only applied to **quantitative** data!
 - For **categorical** data, can't quantify spread since no 'distance' between
 - Instead, give number of categories for given percentile of samples
 - e.g., "90% of samples are in 3 categories"
 - (Pareto chart provides this)


34

Depicting Variation in Charts

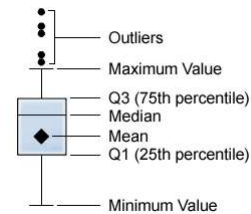
- Histogram (done)
- Cumulative distribution (done)
- Box-and-Whiskers (new)
- Error Bars (new)

35

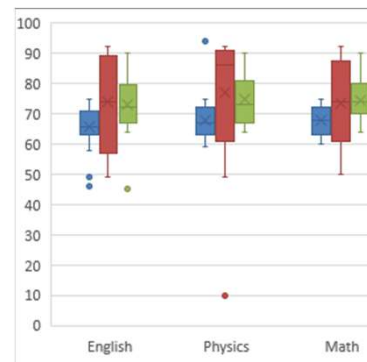
Box-and-Whiskers Chart

- Way of showing variation
- Highlight middle 50% (interquartile range, IQR)
 - “Box”
- Lines go to smallest non-outlier
 - “Whiskers”
- Points indicate outliers
- Middle line shows median
- Sometimes with mean
- **Outlier?** → Data value “way out there”, “far” from the rest
 - Formally, 1.5+ IQRs away from quartile
- Available in Excel 2016 

Sometimes called “boxplot”



<http://support.sas.com/documentation/cdl/en/vaug/65747/HTML/default/images/boxplot.png>



<https://support.office.com/en-us/article/Create-a-box-and-whisker-chart-62f4219f-db4b-4754-ac88-4743f6190f0d>

36

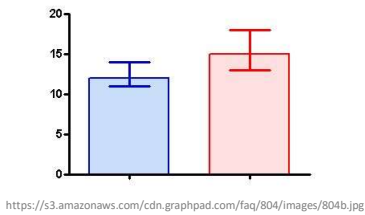
Error Bars

- Line through graph point parallel to axis with "caps"
- Denotes uncertainty (variation) in value
- Often:
 - 1 standard deviation
- Can be (discuss later):
 - 1 standard error
 - 1 confidence interval

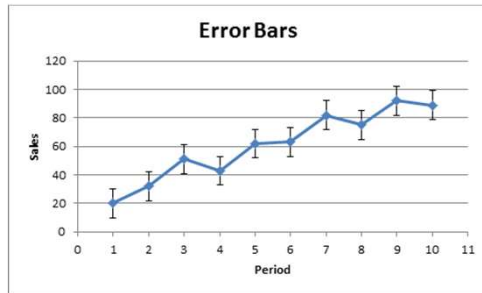
State clearly!



Excel: click "+" → "Error Bars" → "type"



<https://s3.amazonaws.com/cdn.graphpad.com/faq/804/images/804b.jpg>



<http://www.excel-easy.com/examples/images/error-bars/error-bars.png>