Introductory Statistics 4: Data Editing and Summary

Richard Veale

Graduate School of Medicine Kyoto University

https://youtu.be/QAG-gITTXds

Lecture Video at above link



- Types of data (categorical, continuous, etc.)
- Measures of location (mean, median, mode)
- Histograms and Distributions
- Box Plots
- Outliers, Data Trimming
- Measure of spread (variation)
- Scatter Plots

What type of data is this?

Person	Like Ramen?	You a man?
1	Yes	Yes
2	Yes	Yes
3	Yes	No
4	No	Yes
5	No	No
6	No	Yes
7	Yes	Yes
8	Yes	Yes
9	No	Yes
10	Yes	Yes

Are you Man?YesYesNoYes51No31



<u>~-</u>						
len		Yes	No			
an	Yes	5	1			
e E	No	3	1			
. <u>×</u>						

Are you Man?



<u>~-</u>		Are you Man?				
len		Yes	No			
kaπ	Yes	5	1			
е Б	No	3	1			
.≚						

We choose to represent membership in the categories using symbols ("Yes" or "No")

Person	Like Ramen?	You a man?
1	Yes	Yes
2	Yes	Yes
3	Yes	No
4	No	Yes
5	No	No
6	No	Yes
7	Yes	Yes
8	Yes	Yes
9	No	Yes
10	Yes	Yes

Are you Man?YesYesYes5No3

We could use any symbol.... (1 for Yes, 0 for No). This works just the same!

Person	Like Ramen?	You a man?
1	1	1
2	1	1
3	1	0
4	0	1
5	0	0
6	0	1
7	1	1
8	1	1
9	0	1
10	1	1

~-	Are you Man?				
าอเ		1	0		
kaπ	1	5	1		
е Б	0	3	1		

We could use "a" for Yes, "b" for no... We could use 0 for Yes, 1 for No!!! We could name our columns anything we want.

Person	Booboo?	BaaBaa?
1	0	0
2	0	0
3	0	1
4	1	0
5	1	1
6	1	0
7	0	0
8	0	0
9	1	0
10	0	0



We could use "a" for Yes, "b" for no... We could use 0 for Yes, 1 for No!!! We could name our columns anything we want.

Person	Booboo?	BaaBaa?
1	0	0
2	0	0
3	0	1
4	1	0
5	1	1
6	1	0
7	0	0
8	0	0
9	1	0
10	0	0

It is confusing...which is why we use "meaningful" names.

BaaBaa?

		0	1
00	0	5	1
opo	1	3	1
BO			

We could use "a" for Yes, "b" for no... We could use 0 for Yes, 1 for No!!! We could name our columns anything we want.

Person	Booboo?	BaaBaa?
1	0	0
2	0	0
3	0	1
4	1	0
5	1	1
6	1	0
7	0	0
8	0	0
9	1	0
10	0	0

People often think of numbers as representing "amounts"... don't get confused! These are just symbols!

BaaBaa?

		0	1
00	0	5	1
opo	1	3	1
BO			

Bradford - JMP F	Pro					_ 0 <mark>_ X</mark>	ζ
<u>File Edit Tables</u>	<u>R</u> ows <u>C</u> ols	<u>D</u> OE <u>A</u> naly	ze <u>G</u> raph T	<u>o</u> ols <u>V</u> iew	<u>W</u> indow <u>H</u> el	р	
1 🛤 🤮 🗃 🛛	(🗈 🛍 🖶 🦻	à 🔊 📕 🔛	}¤ }k ∎\$ }¤	▫◚ਛੋ₩	🎟 🖷 📕		
■ Bradford	۲ 🔍						_
		BabyGender	Birth Weight	Smoked	Parity		
	1	М	2240	0	0		*
	2	М	3050	0	0		
	3	F	4110	0	2		1
	4	М	3200	1	1		
Columns (4/0)	5	М	3740	0	2		
Birth Weight	6	F	3040	0	1		
Smoked	7	F	2920	0	0		
Parity	8	М	4080	1	0		
_ ,	9	М	2800	0	3		
	10	М	3090	0	2		
	11	М	4110	0	0		
Rows	12	М	3130	0	0		
All rows 100	13	F	2780	0	2		
Selected 0	14	М	2650	0	0		
Excluded 0	15	М	2380	0	0		
Hidden 0	16	М	3170	0	0		
Labelled 0	17	F	2980	0	0		Ŧ
	18	•					

What about this data?

What "types" of data do you see...?

Bradford - JMP Pro								
<u>File Edit Tables</u>	<u>F</u> ile <u>E</u> dit <u>T</u> ables <u>R</u> ows <u>C</u> ols <u>D</u> OE <u>A</u> nalyze <u>G</u> raph T <u>o</u> ols <u>V</u> iew <u>W</u> indow <u>H</u> elp							
: 📴 🦫 💕 💭 🗴 🖦 🖦 🖶 🔊 🖕 : 📴 📴 📴 📾 📾 📾 📾 🔚 🔠 📾 👒 📘								
	Bradford							
		BabyGender	Birth Weight	Smoked	Parity			
		М	2240	0	0			
	2	М	3050	0				
	3	F	4110	0	2		=	
	-	М	3200	1	1			
Columns (4/0)	5	М	3740	0	2			
BabyGender	5	F	3040	0	1			
Smoked	7	F	2920	0	0			
A Parity	3	М	4080	1	0			
		М	2800	0	3			
	1)	М	3090	0	2			
	1.	М	4110	0	0			
Power	12	М	3130	0	0			
All rows 100	18	F	2780	0	2			
Selected 0	1	М	2650	0	0			
Excluded 0	15	М	2380	0	0			
Hidden 0	15	М	3170	0	0			
Labelled 0	17	F	2980	0	0		Ŧ	
	18	•				•		
						☆ 🔲 🔻		

What about this data?

BabyGender

This looks like our "man or not a man".

Probably a category.

Uses "M" and "F" for the symbols for the 2 groups...

Bradford - JMP Pro						
<u>File Edit Tables Rows Cols DOE Analyze Graph Tools View Window Help</u>						
🕮 🍋 💕 🗔 🗶 🛍 🕮 🛤 🔊 🔛 🔛 🔛 📾 📾 🖿 🔚 🖽 📾 🔛 🔛 🔛 🔛						
■Bradford	۲ 🔍					
		BabyGender	Birth Weight	Smoked	Parity	
	1	М	2240	0	0	*
	2	М	3050	0	U	
	3	F	4110	0	2	
	4	М	3200	1	1	
Columns (4/0)	5	М	3740	0	2	
Birth Weight	6	F	3040	0	1	
Smoked	7	F	2920	0	0	
Parity	8	м	4080	1	0	
	9	м	2800	0	3	
	10	м	3090	0	2	
	11	М	4110	0	0	
Rows	12	М	3130	0	0	
All rows 100	13	F	2780	0	2	
Selected 0	14	M	2650	0	0	
Excluded 0	15	M	2380	0	0	
Hidden 0	16	M	3170	0	0	
Labelled 0	17	F	2980	0	0	*
	18	•				

What about this data?

Smoked

This looks similar too.

Did X smoke.

We used "Yes" and "No". It looks like they are using "1" and "0"

(this is common).

Bradford - JMP Pro							
<u>File Edit</u> Tables	<u>R</u> ows <u>C</u> ols	DOE Analy	/ze <u>G</u> raph T <u>o</u>	ols <u>V</u> iew <u>V</u>	<u>W</u> indow <u>H</u> elp)	
🛤 🔁 💕 🗔 🕺 🖻 🛍 🌐 🗷 🔊 🖕 i 🕼 Þ 🛍 🛤 🕽 🧏 🎞 🖬 🗐 i 🕮 🦉							
■Bradford	۲ 🔍						
		BabyGende	Birth Weight	Smoked	Parity		
	1	М	2240	0	0		*
	2	М	3050	0	0		
	3	F	4110	- 0	2		
Columns (4/0)	4	М	3200	1	1		
BabyGondor	5	М	3740	0	2		
Birth Weight	6	F	3040	0	1		
Smoked	7	F	2920	0	0		
A Parity	8	М	4080	1	0		
	9	М	2800	0	3		
	10	M	3090	0	2		
	11	M	4110	0	0		
Rows	12	M	3130	0	0		
All rows 100	13	F	2/80	0	2		
Selected 0	14	м	2650	0	0		
Excluded 0	15	M	2380	0	0		
Hidden 0	10		31/0	0	0		
Labelled 0	1/		2980	0	0		Ψ.
	18			111			

What about this data?

Birth Weight *This looks different!*

In this case, the numbers *actually represent numbers.*

(how much the baby weighed at birth in grams?)

This is data from "Born in Bedford"

This study follows the health of around 13,500 babies born in Bradford between 2007-2010. Did mother smoke or not, etc...



Born in Bradford is one of the biggest and most important medical research studies undertaken in the UK.



"It's like a medical detective story really - trying to piece together the clues in

Nominal data

Categorical data, cannot be put in any specific order. E.g., gender, hair color. **Do you smoke? Yes or No.**

Ordinal data

Categorical data, can be ordered, but the numerical scale is arbitrary (differences between categories unknown).

E.g., pain on a 10-point scale, school grades. Do you smoke (a) none (b) some (c) a lot?

Continuous data

Metric data.

E.g., body temperature, height, air plane speed. How many grams do you smoke?

Discrete data

Also metric, but in integers.

E.g., age in years, number of children. How many cigarettes per day?

In JMP....

Right-click on column, set "data type"

🖳 Bradford - JMP P	ro	C 10 1	* B B 4 5	
File Edit Tables	Rows Cols	DOE Analy	yze Graph Tools View Window	Help
i 🛤 🔁 💕 🗔 🐰	🔁 🛍 🗎 🕖	a 🔊 📙 📴	┡ № # № % ؾ Ħ ⊨ ⊞ Ÿ	
Bradford	<			
		BabyGend	Column Info	1
	1	М		^
	2	М	Standardize Attributes…	0 =
	3	F	Column Properties	2
	4	М	Modeling Type	1
Columns (4/1)	5	М	Droselect Role	2
BabyGender Bitth Weight	6	F	Preselectivole	1
A Smoked	7	F	Formula····	0
A Parity	8	М	Color Cells	0
	9 M Use Value Labels	3		
	10	M		2
	11	M	Label/Unlabel	0
- Paula	12	M	Scroll Lock/Unlock	0
Nows	13	F		2
All rows 100 Selected 0	14	M	Hide/Unnide	0
Evoluded 0	15	M	Exclude/Unexclude	0
Hidden 0	16	M	Data Eiltar	0
Labelled 0	17	F	Data Filter	0 -
	18	•	Sort •	•
			Delete Columns	☆ ■ ▼
			Copy Column Properties	

Protects

In JMP....

 \rightarrow "Data type" refers to the (computer) representation of the data.

 \rightarrow "Modeling type" refers to should it be a category (nominal), or a continuous number, or a ranked ordering?

The Eule Tables	Rows Cols	DOE Analy	ze Graph To	ools View	Window Help)	
) 📑 🔁 🚰 🔜	6 🗈 🛝 🖶 🖉	à 🔊 📕 📴	<u>}</u> ≊ <u>}} ∎\$</u> } ≊	"a ፗ ℍ	- 🛛 🖽 📕		
■Bradford	۲						
		BabyGender	Birth Weight	Smoked	Parity		-
	1	M R	BabyGender -	JMP Pro			
	2		,				
	3 F	F III	'BabyGender' i	n Table 'Bradi	ford'		ОК
Columns (4/1)	4	M	Column Name	BabyGender			Cancel
<mark>⊪ BabyGender</mark> ⊿ Birth Weight	5	F					
	7	F	Data Type Character -		_		Apply
Smoked	8	M			•		Help
Parity	9	M	Modeling Type	Nominal	•		
	10	M					
			Column Prone	rties v			
	11	M	Column Prope	a cico -			
Peur	11 12	M M	Column Prope				
Rows 100	11 12 13	M M F	Column Prope				_
■ Rows All rows 100 Selected 0	11 12 13 14	M F M	2650	0	0		-
Rows All rows 100 Selected 0 Excluded 0	11 12 13 14 15	M F M M	2650 2380	0	0 0		
Rows All rows 100 Selected 0 Excluded 0 Hidden 0	11 12 13 14 15 16	M F M M	2650 2380 3170	0 0 0	0 0 0		
Rows All rows 100 Selected 0 Excluded 0 Hidden 0 Labelled 0	11 12 13 14 15 16 17	M F M M F	2650 2380 3170 2980	0 0 0 0	0 0 0 0		

Measures of Location



"Average":

Arithmetic Mean (of your sample...)

Measures of Location



Doesn't make sense for categorical (or ordinal) data.

→ What is the average of "red hair", "blond hair" and "brown hair"?

 \rightarrow What is the average of "none" and "sometimes"?

Median

Median

1) Sort the n observations from smallest to largest

2a) The median is the single middle value if n is odd.2b) The median is the average of the two middle values if n is even.

Mode

Mode

Most frequent value in observations.

E.g.:

Data: 1123555668

Mode = 5

Can be used for categorical, ordinal, and metric (continuous, discrete) data.

Visualizing Data: Histogram



Visualizing Data: Histogram



Household income in Japan (www.stat.go.jp, 2004)

In JMP

🖼 Bradford - JMP Pro							
File Edit Tables Re	ows Cols	DOE	Ana	lyze Graph Tools V	/iew	Wind	dow Help
🗄 🎦 🚰 🗔 🔏 🗈	n 🛍 🖶 🖉	à 🔊 🖕	F	Distribution			Distribution of a batch of values.
■Bradford D			<u>у</u> х	Fit Y by X			Frequencies if categorical. Means
		BabyG	¥	Matched Pairs			Histograms, Box Plots, Quantile
	1	M		Tabulate			- Plots. Tests on means, Fitting
	3	F					distributions. Capability.
	4	М	>	Fit Model			1
Columns (4/1)	5	М		Modeling		•	2
Birth Weight	6	F		Multivariate Methods		•	1
Smoked	/ 8	F M		Quality and Process			0
A Parity	9	M		Poliability and Suprival			3
	10	М				·	2
	11	М		Consumer Research		•	0
Rows	12	M	_	3130	0		0
All rows 100	13	г М		2/80	0		0
Selected 0	15	M		2380	0		0
Hidden 0	16	М		3170	0		0
Labelled 0	17	F		2980	0		0 -
	18	•					

In JMP



"Distribution" in JMP

🔚 Bradford - Distribution of Birt... 💶 💷 💻 💴



Shows histogram + Boxplot

Shows quantiles

Presents summary statistics

What is a "quartile"?

1st quartile: 25% of data below, 75% above

- 1) Sort the n observations from smallest to largest value.
- 2) The value that separates the data into 25% below and 75% above is the 1st quartile.



Quartiles useful for: *Outliers*

1st quartile: 25% of data below, 75% above 3rd quartile: 75% of data below, 25% above

Inter-quartile range (IQR) = 3rd quartile - 1st quartile



Box plot by variable in JMP

Menu-> Graph -> Graph builder

X: Smoked Y: Baby weight



Spread of data (variance)

Population Variance ("Sigma Squared")

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

Sample Variance

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \hat{x})^{2}$$

Spread (variance)

 $s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \hat{x})^{2}$

How far is each individual data point from the average? Squared to make positive and over-weigh data further from average.



Variance versus Standard Deviation

Sample Variance:

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \hat{x})^{2}$$

Sample Standard Deviation "s"

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

Statistic versus Parameter

- I keep saying
- → "Sample standard deviation"
- → "Sample average"

Why do I say "sample"? What is so important?

Statistic versus Parameter

We differentiate between:

Statistic : We actually measure this. Using a sample of the population.

Parameter : Imaginary, (unknowable) theoretical property of the true population.

Statistic versus Parameter

We differentiate between:

Statistic : We actually measure this. Using a sample of the population.

Parameter : Imaginary, (unknowable) theoretical property of the true population.

The key theory is that we can *estimate* parameters using *samples*. And use it!

Sample Average is *Unbiased Estimator* of the Population Mean

No matter our sample size, the sample average will *underestimate* and *overestimate* the true population mean an equal amount!

 \rightarrow In other words, the average of our sample averages will equal the true population mean...

Sample Average is *Unbiased Estimator* of the Population Mean

No matter our sample size, the sample average will *underestimate* and *overestimate* the true population mean an equal amount!

 \rightarrow In other words, the average of our sample averages will equal the true population mean...

...meta.

Biased Estimator of Variance...

We take 1/(n-1) for sample variance.

But for population variance (parameter) it is 1/n.

Why?



Biased Estimator of Variance...

This is because unlike the mean, variance is a **biased** estimator.

Using *n*-1 rather than *n* is called Bessel's Correction – it is one method of correcting sample variance.

Unfortunately, square rooting the variance to get standard deviation "uncorrects" it some...oh well.

Consider these numbers as the values of the whole population:



If we take a small sample (n=3), we have a good chance that we estimate μ and σ well with sample mean (x[^]) and sample variance (s):

But sometimes, we will randomly pick sample elements that do not cover μ and in this case it is obvious that s will be much smaller than σ :

Consider these numbers as the values of the whole population:



If we take a small sample (n=3), we have a good chance that we estimate μ and σ well with sample mean (x[^]) and sample variance (s):



But sometimes, we will randomly pick sample elements that do not cover μ and in this case it is obvious that s will be much smaller than σ :











Doing n-1 is a "hack" (rough correction) to get the values to underestimate the variance less in the limit.

It works because formally sample variance has one less degree of freedom (if we know the mean, and we know all but one sample, we can cheat and calculate the value of the last sample! Giving us information out of nowhere?

So, what that means is basically that we have less information than we thought.

If we knew the true mean μ and used it \rightarrow unbiased.

See it work

Let's say, we have 1,000,000 birth weight values [g] μ =3200 g, σ =560 g (simulated with normal distribution)



See it work

Mean of SD estimated for 100 samples with sample sizes 2 to 10.



51

See it work



Multiply by n/(n-1)



Bar Graph in JMP

Menu \rightarrow Graph \rightarrow Graph builder

X: Baby gender; Y: Baby weight



Scatter Plot in JMP

Menu-> Graph -> Graph builder

DE Analyze Graph Tools View Windo	w Help	
▖▁ੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵੵ	Par 🔒 💡	
Bradford_BabyMotherWeight - Graph B	uilder - J	MP Pro
Graph Builder	•	
Recall Start Over Done	\sim	
Variables		Title
Columns		litte
BabyWeight		Group X
⊿ Points		
Jitter 🔽		
Statistic None		
Error Bars None -		
Variables		
	Y	Drag variables into drop zones
	Map Shape	Х
	strape	

Scatter plots can visualize correlations between continuous variables

Scatter Plot in JMP

56

X: Mother weight; Y: Baby weight

