# STRUCTURING YOUR DATA THE RIGHT WAY THE FIRST TIME

Sandhya Kambhampati @sandhya\_\_k



## PDFS HBAAAAA

ARE THE WORST

#### GOVERNMENT FILE STRUCTURE

Name of the last	С	D	E	F	G					
			Sonderausv	vertung						
1	Personal am 15.12.2013									
9			in Pflegeheimen auf Kr	reisebene in Bayern						
				Stationä	re Pflege					
Überwiegender Tätigkeitsbereich im Pflegeheim										
	Damanal Income	Auszubildende/-r,								

Personal Insgesamt zusätzliche Betreuung (Um-)Schüler/-in Pflege und Betreuung soziale Betreuung (§ 87b SGB XI) 917 59 599 38 32 5 2 1 0 443 3 675 227 171

### FOIA DOCUMENTS

Α	В	С	D	E	F	G	Н	1	J	K
l	Ticket Sales.	\$103,410.68	fa sa ex	culty and students, a les for conference ar	nd nd	I for sales of admissi I money received for national tournament value paid by ticket pass).	shipping and s that are pas	d handling ss-through	of tickets. Do national research	ot include ticket eport amounts in
		Men's Teams Only	П	Women's Teams Only		Not Allocated by Gender		Summary	Men's	\$97,432.97
	Revenues by Source	Ticket Sales.	П	Ticket Sales.	П	Ticket Sales.			Women's	\$5,977.71
		1	П	1	П	1			No Gender	\$0.00
	Archery		П		П				Grand Total	\$103,410.68
	Badmiton		11		П					
	Baseball		1 1		П					

## THE BASICS

- Rows
- Columns
- Column Headers
- Cells

COLUMN

COLUMN HEADER

**ROW** 

	A	В	С	D
	first_name	last_name	phone_number	email_address
2				
3				
1				
5				
5				
7				
3				
0 (		)		
	CELL			

## WHY DOES STRUCTURE MATTER?

### ARE THESE THE SAME?

Α	В	С		
Student	test_score_1	test_score_2		
<b>Bob Smith</b>	88	67		
Anita Doe	100	100		
John Phillips	76	87		

Student	Bob Smith	Anita Doe	John Phillips
test_score_1	88	100	76
test_score_2	67	100	87

## THINK ABOUT FUNCTIONS

SORT
AGGREGATE
TRANSFORM
FILTER

## STOP THIS!





	Α	В	С	D	E	F
1	unique id	total3	total2	namef	namel	other1
2	1	64328	392902	Sam	Smith	738393
3	2	42910	2048349	Bob	Miller	28102
4	3	27293	2302938	Sandhya	Kambhampa	1281912
5	4	91213	2392302	Sam	Jones	12182913
6	5	42323	23826	Nikil	Patel	4289328943
7	6	492303	3422	Samantha	Smith	32849
8	7	23900	3423422	John	Doe	29401

## EXAMPLE FILE

http://bit.ly/2lw2JAV

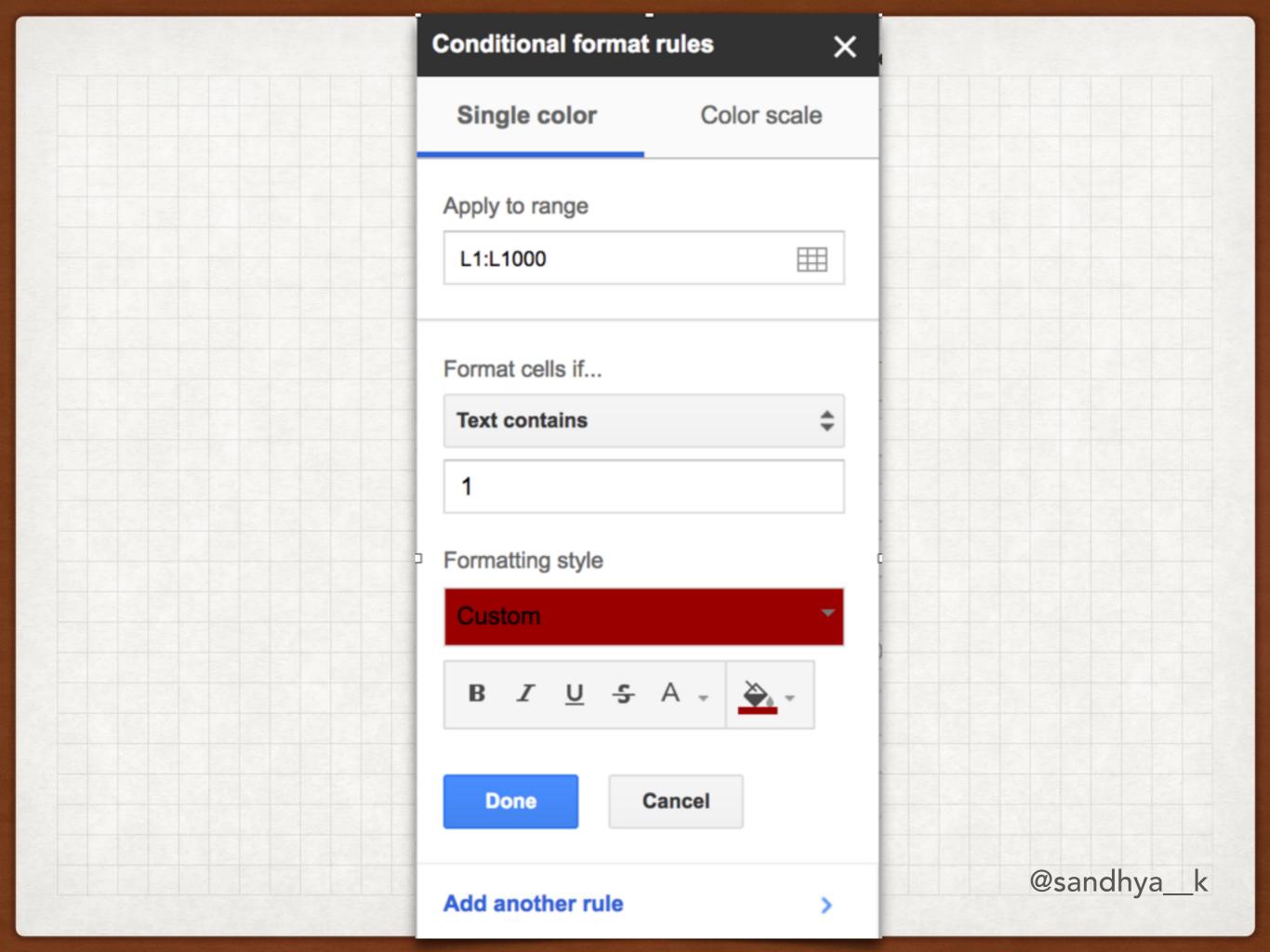
#### INSTEAD OF THIS

Α	В	C	D	E	F	G
ID	Player	Date_of_birth	Day	Month	Year	Age
1	Sergio Romero	2/22/87	22	February	1987	29
2	Nahuel Guzmán	2/10/86	10	February	1986	30
3	Mariano Andújar	7/30/83	30	July	1983	33
4	Pablo Zabaleta	1/16/85	16	January	1985	31
5	Martín Demichelis	12/20/80	20	December	1980	35
6	Marcos Rojo	3/20/90	20	March	1990	26
7	Nicolás Otamendi	2/12/88	12	February	1988	28
8	Ramiro Funes Mori	3/5/91	5	March	1991	25
9	Facundo Roncaglia	2/10/87	10	February	1987	29
10	Gabriel Mercado	3/18/87	18	March	1987	29
11	Emmanuel Más	1/15/89	15	January	1989	27
12	Mateo Musacchio	8/26/90	26	August	1990	25
13	Javier Mascherano	6/8/84	8	June	1984	32
14	Ángel Di María	2/14/88	14	February	1988	28

### USE FLAGS

Α	В	С	D	Е	F	G	Н
ID	Player	Date_of_birth	Month	Day	Year	Age	Checked
1	Sergio Romero	2/22/87	February	22	1987	30	1
2	Nahuel Guzmán	2/10/86	February	10	1986	31	1
3	Mariano Andújar	7/30/83	July	30	1983	33	0
4	Pablo Zabaleta	1/16/85	January	16	1985	32	1
5	Martín Demichelis	12/20/80	December	20	1980	36	1
6	Marcos Rojo	3/20/90	March	20	1990	26	1
7	Nicolás Otamendi	2/12/88	February	12	1988	29	0
8	Ramiro Funes Mori	3/5/91	March	5	1991	25	0
9	Facundo Roncaglia	2/10/87	February	10	1987	30	0
10	Gabriel Mercado	3/18/87	March	18	1987	29	0
11	Emmanuel Más	1/15/89	January	15	1989	28	0
12	Mateo Musacchio	8/26/90	August	26	1990	26	0

File Edit View	Insert	Format Data Tools Add-o	ons He	lp	All changes save	ed in Drive		
	\$ %	Number	<b>&gt;</b>	¥	B I 5	<u>A</u> - <u>\</u>	<b>→</b> ⊞ →	1 ,
		Font size	<b>&gt;</b>					
В		1 011 0120			F	G	Н	
layer	Date_	B Bold	₩B		Year	Age	Weekday	Cap
ergio Romero		I Italic	ЖI	22	1987	30	1	
ahuel Guzmán		<u>U</u> Underline	ЖU	10	1986	31	2	
lariano Andújar		Strikethrough Option+	Shift+5	30	1983	33	7	'
ablo Zabaleta		•		16	1985	32	4	
lartín Demichelis		Align	<b>&gt;</b>	20	1980	36	7	'
larcos Rojo		Merge cells	<b>&gt;</b>	20	1990	26	3	3
icolás Otamendi		Text wrapping	<b>&gt;</b>	12	1988	29	6	5
amiro Funes Mori		Text rotation	<b>&gt;</b>	5	1991	25	3	1
acundo Roncaglia				10	1987	30	3	1
abriel Mercado		Conditional formatting		18	1987	29	4	
mmanuel Más		Alternating colors		15	1989	28	1	
lateo Musacchio				26	1990	26	1	
vier Mascherano		$\mathcal{I}_{X}$ Clear formatting	#\	8	1984	32	6	i
ngel Di María		2/14/88 February		14	1988	29	1	
D		C /20 /00 I		20	1000	20		



#### ranks\_pflege\_datadictionary.md ...

#### ## Data Dictionary for `ranks\_pflege\_data.xlsx`

\*\*Source\*\*: Data comes from Statisches Landesamt for December 2013.

\*\*Note\*\*: This file contains the rankings for each county for the major data questions we asked.

#### ### Contents

\*All files contain the same two beginning A and B columns (`Name` and `State`), where `name` means the name of the district or county and `state` is the state the county is found in.\*

#### #### Sheets

#### \* \*\*More than 1 bed\*\*

- `Rank`: ranking of the total # beds that are in rooms with more than 1 beds in the county or district in descending order
- `Total of >1bed`: total number of beds that are in rooms with more than 1 bed
- \*\*About this data\*\*: the rankings of the data show that `Berlin` has the highest number of beds overall at `11,896`. This shows that `Berlin` has a lot more older nursing homes than other places as they have a lot more rooms with more than 1 bed.

#### \* \*\*Ratio of beds\*\*

- `Total of >1 beds`: total number of beds that are in rooms with more than 1 bed
- `All beds`: total number of beds in the county or district
- `Ratio of >1bed to allbeds`: `total of >1 beds`/ `all beds`; in other words, what percentage of beds are in rooms with more than 1 bed
- 'Rank': rank of the 'ratio of >1bed to allbeds'

#### DATA DICTIONARY

#### GOOD HEADERS & CLEAN DATA

4	Α	В	С	D	E	F
1	unique_id	total_debt	total_revenue	first_name	last_name	taxes
2	1	64328	392902	Sam	Smith	738393
3	2	42910	2048349	Bob	Miller	28102
4	3	27293	2302938	Sandhya	Kambhampa	1281912
5	4	91213	2392302	Sam	Jones	12182913
6	5	42323	23826	Nikil	Patel	4289328943
7	6	492303	3422	Samantha	Smith	32849
8	7	23900	3423422	John	Doe	29401

## QUESTIONS? THANK YOU!