

spreadsheets



 [@JennyBryan](https://twitter.com/JennyBryan)

 [@jennybc](https://github.com/jennybc)

STAT
545  [@STAT545](https://twitter.com/STAT545)
 <http://stat545.com>

relevant links, credits, and slides:

https://github.com/jennybc/2016-06_spreadsheets

Rich FitzJohn



Research Software Engineer
University College London

 [@rgfitzjohn](https://twitter.com/rgfitzjohn)

 [@richfitz](https://github.com/richfitz)



spreadsheets: a dystopian moonscape of unrecorded user actions

— Gordon Shotwell

some of my best
friends use
spreadsheets

I supported myself for
~4 years doing
spreadsheets

~1 billion use Microsoft Office

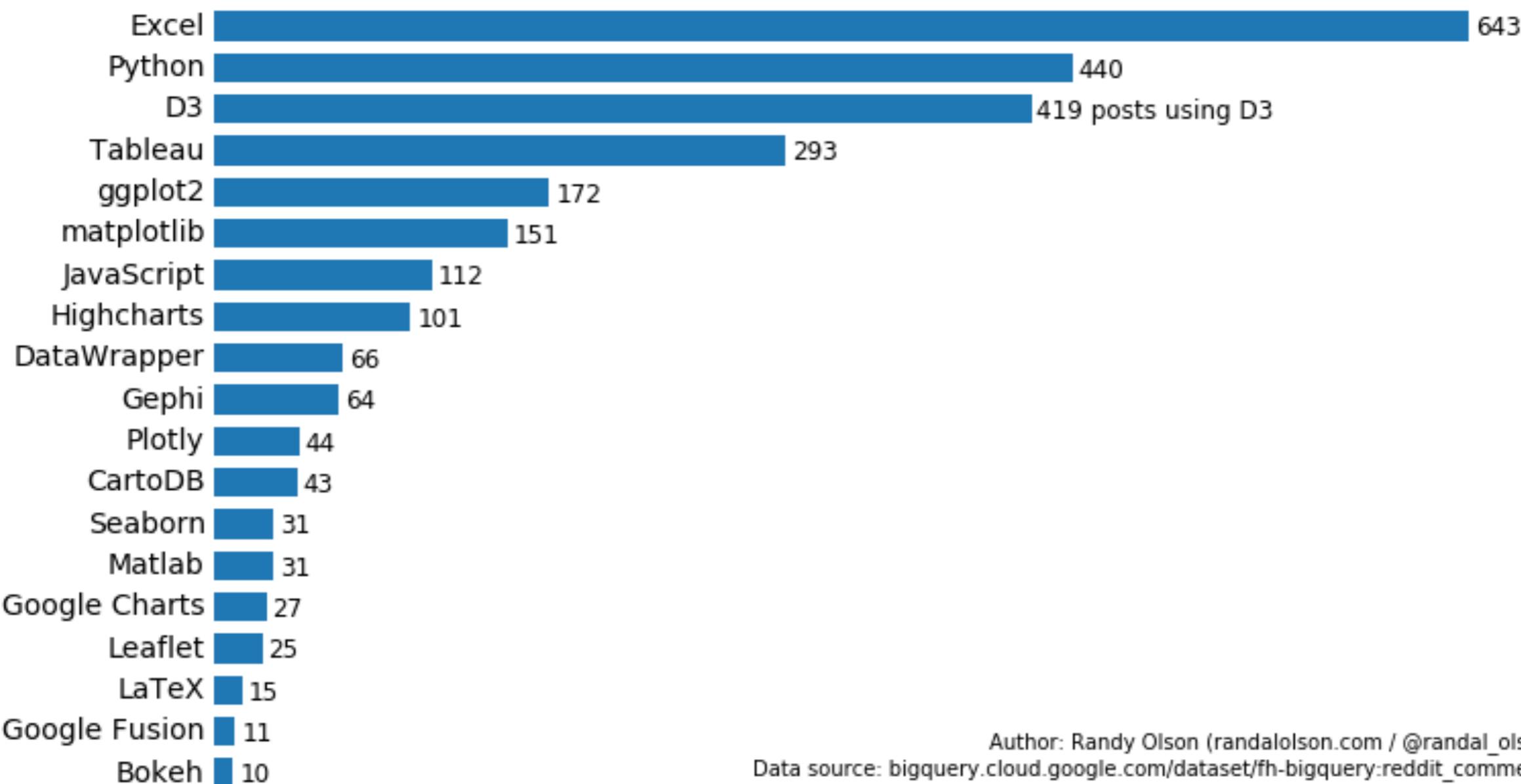
~650 million use spreadsheets

>50% use formulas

1 - 5 million people use Python

250K - 1 million people use R

Tools used in /r/DataIsBeautiful OC posts, 2014-2016



Author: Randy Olson (randalolson.com / [@randal_olson](https://twitter.com/randal_olson))

Data source: bigquery.cloud.google.com/dataset/fh-bigquery:reddit_comments

you go into data analysis
with the tools you know,
not the tools you need

spreadsheets combine:
data
logic
figures
formatted tables
+ reactivity

spreadsheets users
use workbooks
like I would use
a data analytic git repo

a data analytic project:

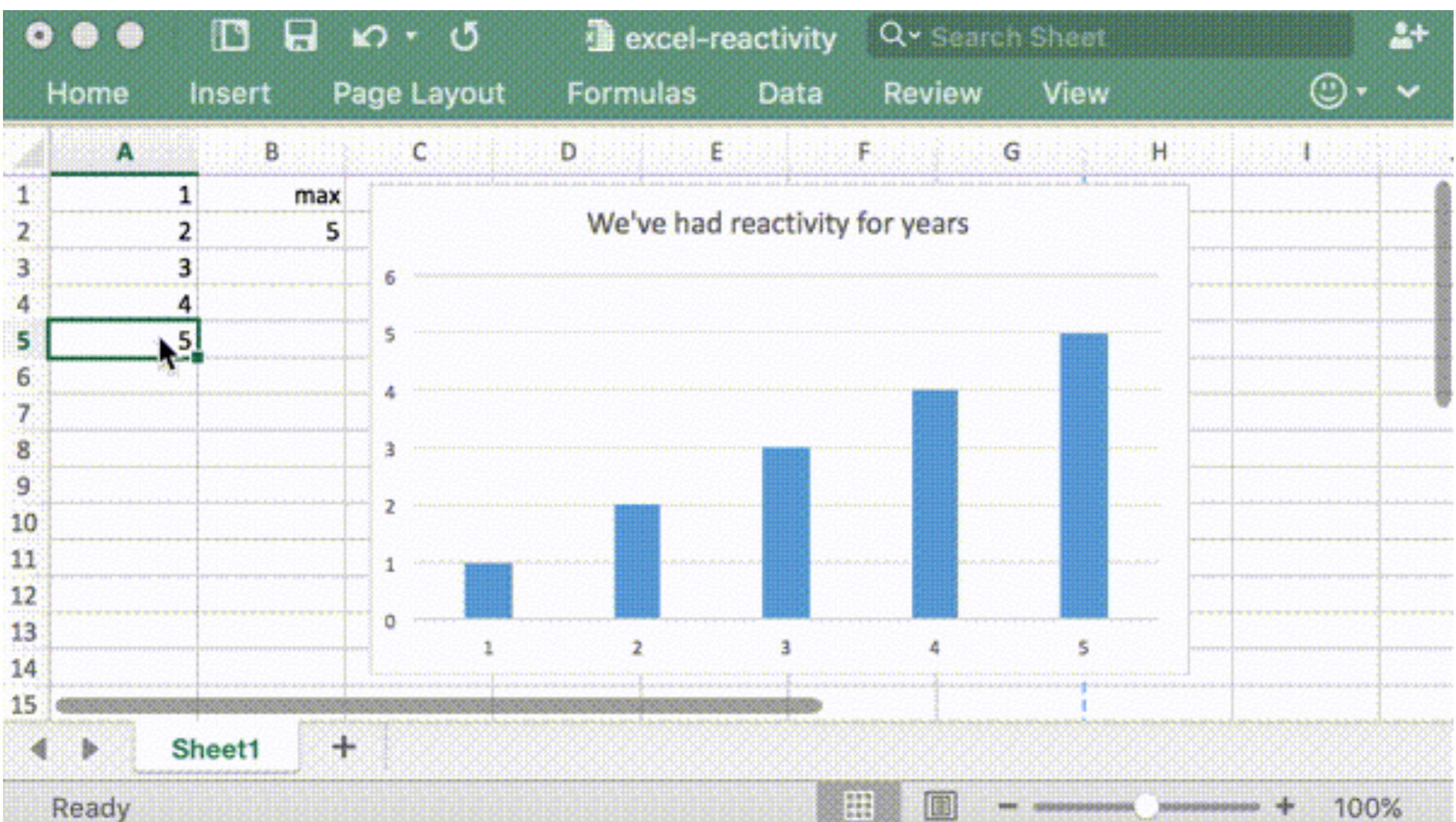
data

.R, .Rmd

.png, .svg

.md, .html, .pdf, Shiny app

+ build and deploy





**syntax
bullshittery**

spreadsheets are not
going away

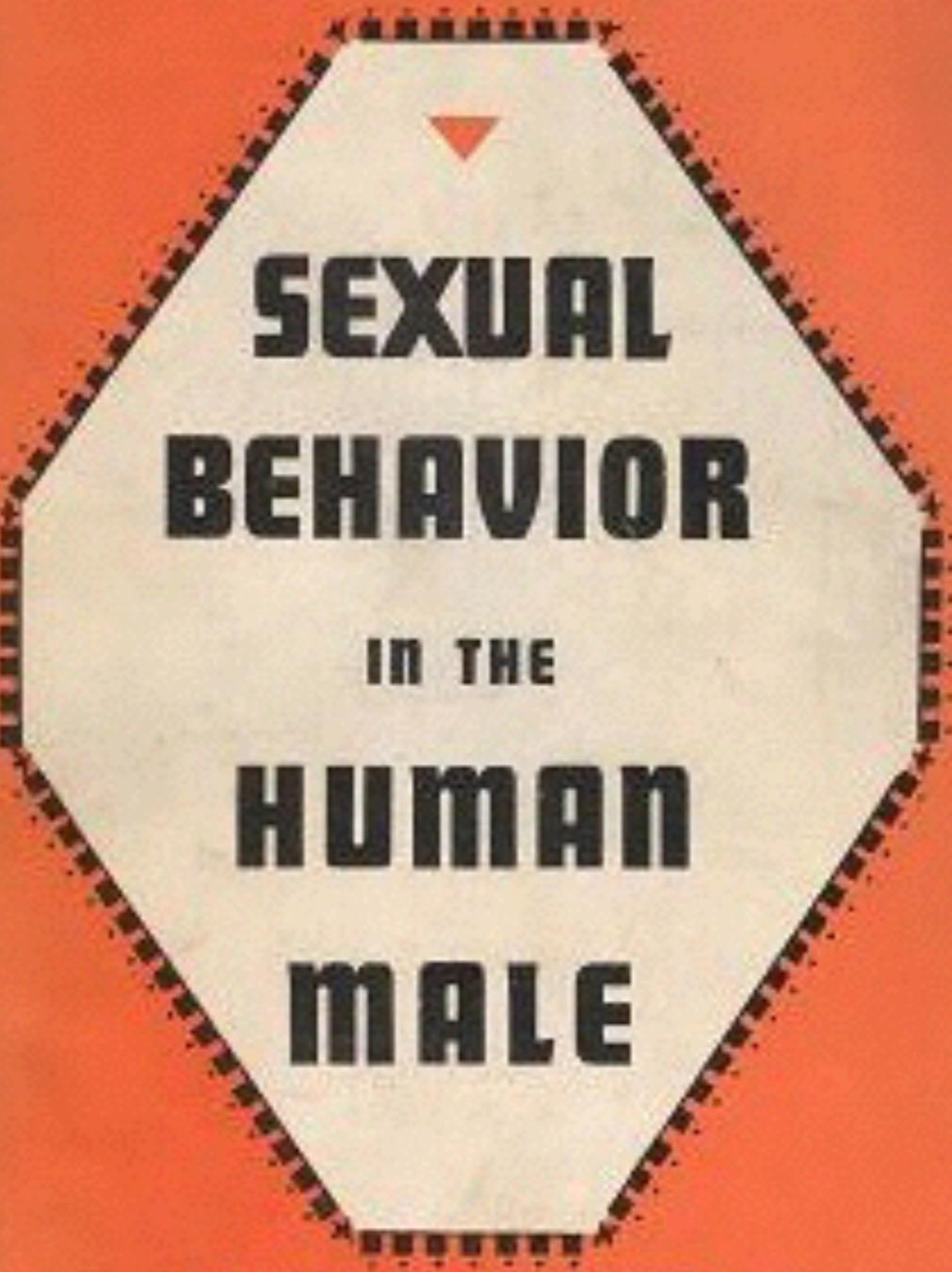
deal with it



Jenny Bryan @JennyBryan · Apr 20

I'm seeking TRUE, crazy spreadsheet stories. Happy to get the actual sheet or just a description of the crazy. Also: I can keep a secret.

Based on surveys made by members of the Staff of Indiana
University and supported by the National Research Council
with Rockefeller Foundation funds



**SEXUAL
BEHAVIOR
IN THE
HUMAN
MALE**

ALFRED C. KINSEY
WARDELL B. POMEROY
CLYDE E. MARTIN

Based on surveys made by members of the Staff of Indiana
University and supported by the National Research Council
with Rockefeller Foundation funds

SPREADSHEET
BEHAVIOR
IN THE
HUMAN
MALE

ALFRED C. KINSEY

WARDELL B. POMEROY

CLYDE E. MARTIN

Based on surveys made by members of the Staff of Indiana University and supported by the National Research Council with Rockefeller Foundation funds

SPREADSHEET
BEHAVIOR
IN THE
HUMAN
MALE

ALFRED C. KINSEY
WARDELL B. POMEROY
CLYDE E. MARTIN

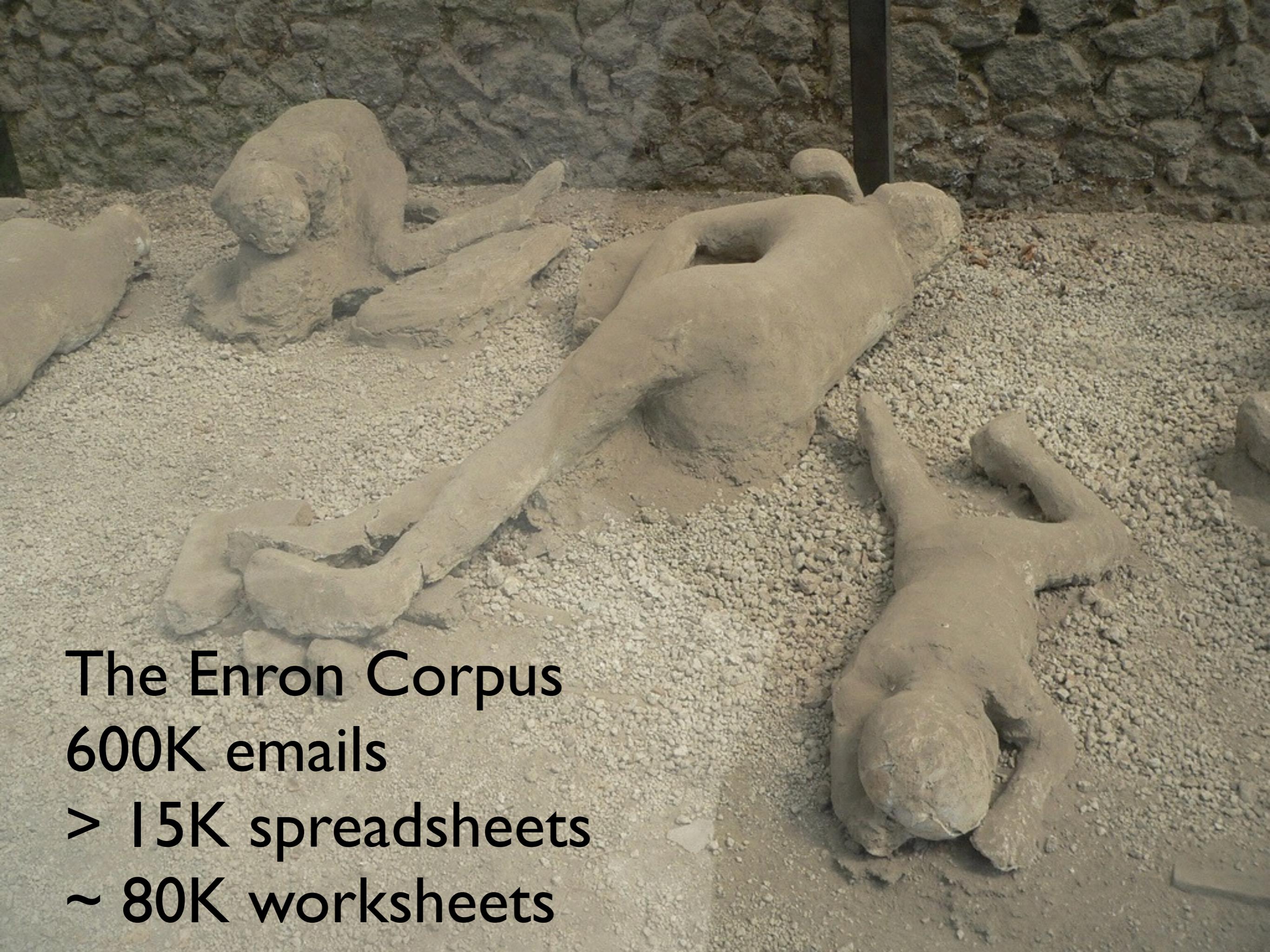
what you **THINK** people are doing

!=

what you think people **SHOULD** be doing

!=

what people **ARE ACTUALLY** doing



The Enron Corpus
600K emails
> 15K spreadsheets
~ 80K worksheets



Enron North America - West Gas

November 9, 2001



ENA - West Gas Contacts

Houston Office

Barry Tycholiz	(713) 853-1587
Kim Ward	(713) 853-0685
Stephanie Miller	(713) 853-1688
Philip Polsky	(713) 853-5181

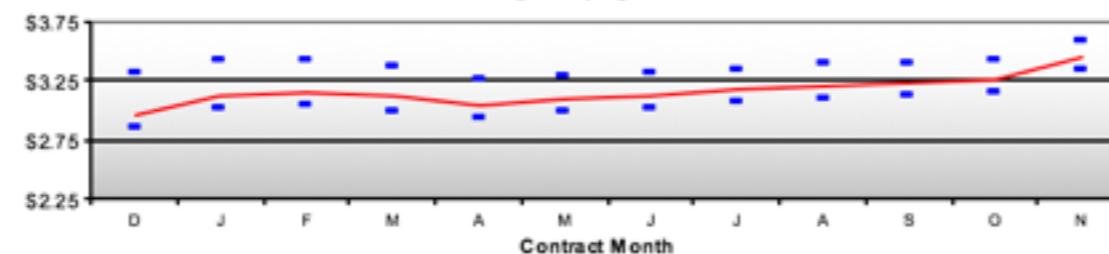
Regional Offices

Mark Whitt	(303) 575-6473	Denver
Paul Lucci	(303) 575-6474	Denver
Tyrell Harrison	(303) 575-6478	Denver
Dave Fuller	(503) 464-3732	Portland

Forward Prices (US\$/MMBtu)

NYMEX	
	SETTLE
Cash	
ROM	
Dec-01	2.960 0.090
Dec-01 to Mar-02	3.088 0.083
Apr-02 to Oct-02	3.166 0.084
Nov-02 to Mar-03	3.651 0.090
One Year Strip*	3.165 0.084

Forward NYMEX Strip
with trailing 10-day highs/lows



IF NWPL Rocky Mountains

Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.890	1.910		
2.060	2.080		
2.395	2.415	(0.565)	(0.545)
2.594	2.614	(0.494)	(0.474)
2.581	2.601	(0.585)	(0.565)
3.356	3.376	(0.295)	(0.275)
2.634	2.654	(0.530)	(0.510)

IF CIG Rocky Mountains

Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.940	1.960		
1.960	1.980		
2.345	2.365	(0.615)	(0.595)
2.548	2.568	(0.540)	(0.520)
2.471	2.491	(0.695)	(0.675)
3.311	3.331	(0.340)	(0.320)
2.551	2.571	(0.614)	(0.594)

IF EL Paso Permian

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.375	2.395		
2.420	2.440		
2.700	2.720	(0.260)	(0.240)
2.855	2.875	(0.233)	(0.213)
3.009	3.029	(0.158)	(0.138)
3.499	3.519	(0.153)	(0.133)
2.982	3.002	(0.182)	(0.162)

IF EL Paso San Juan

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.450	2.470		
2.350	2.370		
2.560	2.580	(0.400)	(0.380)
2.743	2.763	(0.345)	(0.325)
2.801	2.821	(0.365)	(0.345)
3.421	3.441	(0.230)	(0.210)
2.817	2.837	(0.347)	(0.327)

AECO / NIT

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.376	2.396		
2.398	2.418		
2.552	2.572	(0.408)	(0.388)
2.616	2.636	(0.472)	(0.452)
2.661	2.681	(0.505)	(0.485)
3.216	3.236	(0.435)	(0.415)
2.676	2.696	(0.488)	(0.468)

IF NWPL Canadian Border (Sumas)

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.480	2.500		
2.460	2.480		
2.800	2.820	(0.160)	(0.140)
2.892	2.912	(0.196)	(0.176)
2.796	2.816	(0.370)	(0.350)
3.706	3.726	0.055	0.075
2.880	2.900	(0.285)	(0.265)

IF PEPL TX-OK

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.530	2.550		
2.530	2.550		
2.828	2.848	(0.133)	(0.113)
2.958	2.978	(0.130)	(0.110)
3.046	3.066	(0.120)	(0.100)
3.531	3.551	(0.120)	(0.100)
3.041	3.061	(0.123)	(0.103)

	EUSES	Enron
Number of spreadsheets analyzed	4,447	15,770
Number of spreadsheets with formulas	1,961	9,120
Number of worksheets	16,853	79,983
Maximum number of worksheets	106	175
Number of non-empty cells	8,209,095	97,636,511
Average number of non-empty cells per spreadsheet	1,846	6,191
Number of formulas	730,186	20,277,835
Average of formulas per spreadsheet with formulas	372	2,223
Number of unique formulas	65,143	913,472
Number of unique formulas per spreadsheet with formulas	33	100

from Hermans, Murphy-Hill

data in formatting

Plot: 2			
Date collected	Species	Sex	Weight
1/8/14	NA		
1/8/14	DM	M	44
1/8/14	DM	M	38
1/8/14	OL		
1/8/14	PE	M	22
1/8/14	DM	M	38
1/8/14	DM	M	48
1/8/14	DM	M	43
1/8/14	DM	F	35
1/8/14	DM	M	43
1/8/14	DM	F	37
1/8/14	PF	F	7
1/8/14	DM	M	45
1/8/14	OT		
1/8/14	DS	M	157
1/8/14	OX		
2/18/14	NA	M	218
2/18/14	PF	F	7
2/18/14	DM	M	52

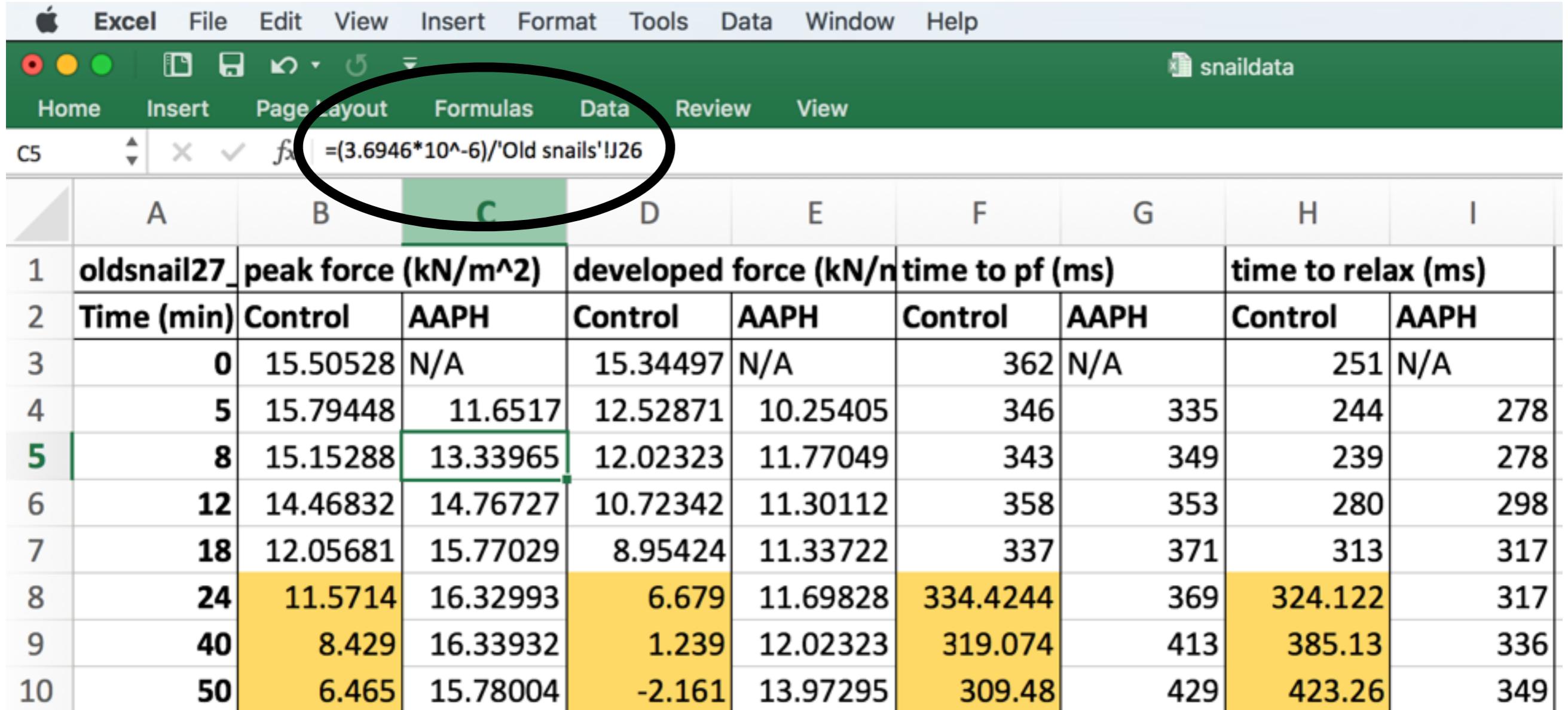
measurement device not calibrated

A	B	C
1 species		tail length
2 Allactaga balikunica		177.32
3 Allactaga bullata		165.2
4 Allocricetulus eversmanni		18.64
5 Apodemus uralensis		84.89
6 Arvicola amphibius		105.14
7 Brachytarsomys albicauda		230.02
8 Brachyuromys betsileoensis		83.32
9 Cardiocranus paradoxus		74.36
10 Castor fiber		379.89
11 Cricetulus barabensis		24.72
12		
13 bold = needs checking		
14 yellow = from different source		
15		

small multiples

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	
oldemail	peak force (kN/m ²)	developed force (kN)	Time to p (ms)	Time to relax (ms)																													
Time [m]	Control	AAPH	Control	AAPH	Control	AAPH	Control	AAPH																									
0	15.51	n/a	15.34	n/a	362	n/a	251	n/a																									
5	15.79	11.85	12.53	10.25	346	335	244	278																									
8	15.15	13.34	12.01	11.77	343	349	239	278																									
12	14.47	14.77	10.72	11.3	358	353	280	298																									
18	12.06	15.77	8.554	11.34	337	371	313	317																									
24	11.57	16.33	6.679	11.7	334.4	369	324.1	317																									
40	8.429	16.34	1.235	12.02	319.1	413	385.1	336																									
50	6.465	15.78	-2.163	13.57	309.5	429	423.3	349																									
oldemail	peak force (kN/m ²)	developed force (kN)	Time to p (ms)	Time to relax (ms)																													
Time [m]	Control	AAPH	Control	AAPH	Control	AAPH	Control	AAPH																									
0	9.373	n/a	8.385	n/a	429	n/a	545	n/a																									
5	6.764	3.239	4.924	2.549	357	341	569	558																									
8	5.66	3.446	4.318	2.772	364	353	497	563																									
12	6.394	5.176	3.838	3.571	362	349	656	570																									
18	5.19	5.553	2.558	3.678	352	382	789	622																									
24	3.565	5.048	0.351	3.331	316	360	826.2	585																									
40	0.333	4.649	-4.477	3.598	261.2	381	1050	594																									
50	-1.687	5.902	-7.494	4.211	227	392	1189	581																									
oldemail	peak force (kN/m ²)	developed force (kN)	Time to p (ms)	Time to relax (ms)																													
Time [m]	Control	AAPH	Control	AAPH	Control	AAPH	Control	AAPH																									
0	6.546	n/a	6.842	n/a	440	n/a	459	n/a																									
5	7.551	4.360	6.427	3.835	361	365	301	393																									
8	6.886	4.514	6.012	3.574	381	379	314	403																									
12	5.763	5.123	4.590	3.539	374	369	370	426																									
18	4.329	4.357	3.325	3.835	357	340	424	491																									
24	3.768	4.385	2.236	3.45	326.3	330	375.6	480																									
40	1.141	3.886	-1.089	3.415	267.7	350	377.7	465																									
50	-0.500	4.043	-3.167	3.443	231.2	350	379	436																									
oldemail	peak force (kN/m ²)	developed force (kN)	Time to p (ms)	Time to relax (ms)																													
Time [m]	Control	AAPH	Control	AAPH	Control	AAPH	Control	AAPH																									
0	21.82	n/a	20.78	n/a	325	n/a	321	n/a																									
5	16.35	14.78	15.35	11.55	307	348	331	338																									
8	15.77	15.7	13.48	12.18	313	344	334	330																									
12	13.42	15.77	9.345	9.741	343	365	452	436																									
18	11.12	13.51	7.325	8.269	359	366	491	624																									
24	6.295	11.79	2.205	7.023	366.7	346	550.3	644																									
40	-2.294	12.81	-9.401	8.722	405.4	385	721.2	531																									
50	-8.062	12.33	-16.66	7.079	429.5	421	828	689																									
oldemail	peak force (kN/m ²)	developed force (kN)	Time to p (ms)	Time to relax (ms)																													
Time [m]	Control	AAPH	Control	AAPH	Control	AAPH	Control	AAPH																									
0	6.573	n/a	6.254	n/a	532	n/a	627	n/a																									
5	7.446	3.778	6.153	2.673	510	539	385	612																									
8	6.475	3.591	5.372	2.565	516	545	393	644																									
12	6.542	3.82	4.357	2.332	532	512	426	628																									
18	5.352	3.563	3.595	2.156	573	518																											

data in formulas



The screenshot shows a Microsoft Excel spreadsheet titled "snaildata". The ribbon menu is visible at the top, with the "Formulas" tab highlighted. The formula bar shows the formula $=(3.6946*10^{-6})/\text{'Old snails'!J26}$. The main table has columns labeled A through I. The first row contains column headers: "oldsnail27", "peak force (kN/m^2)", "developed force (kN/n", "time to pf (ms)", and "time to relax (ms)". The second row contains headers for "Time (min)" and "Control" and "AAPH" for each of the four columns. The data rows show measurements for times 0, 5, 8, 12, 18, 24, 40, and 50 minutes. Cells C5 and C8 are highlighted with green borders. Cells C5, C8, D8, E8, F8, G8, H8, and I8 are circled in black.

oldsnail27	peak force (kN/m ²)	developed force (kN/n	time to pf (ms)	time to relax (ms)
Time (min)	Control	AAPH	Control	AAPH
0	15.50528	N/A	15.34497	N/A
5	15.79448	11.6517	12.52871	10.25405
8	15.15288	13.33965	12.02323	11.77049
12	14.46832	14.76727	10.72342	11.30112
18	12.05681	15.77029	8.95424	11.33722
24	11.5714	16.32993	6.679	11.69828
40	8.429	16.33932	1.239	12.02323
50	6.465	15.78004	-2.161	13.97295

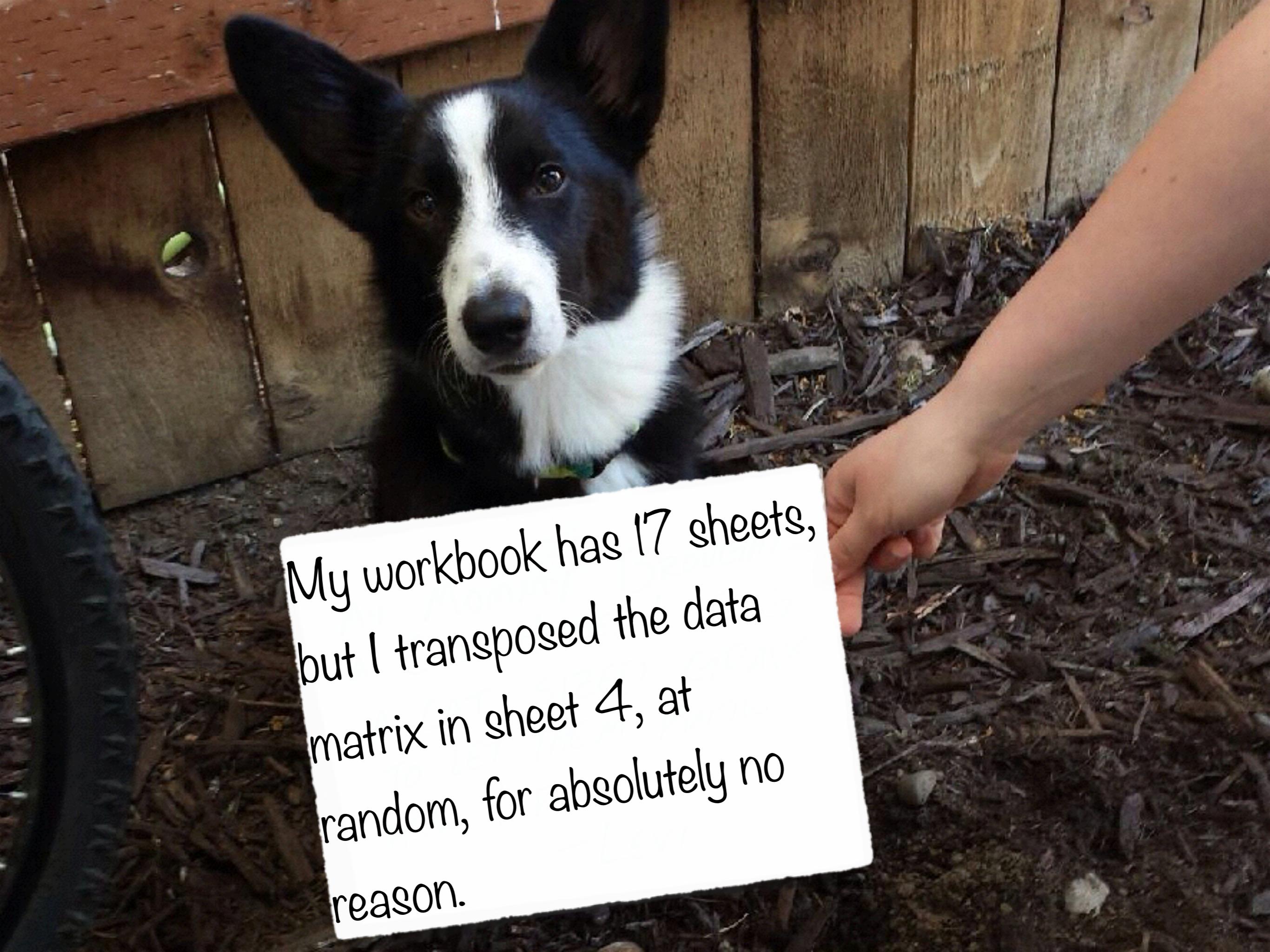
$=(3.6946*10^{-6})/\text{'Old snails'!J26}$

data in (merged) column headers

Snake River wild chinook master

	A	B	C	D	E	F	G	H	I	J	K	L											
1	Estimated wild spring and summer chinook salmon based on clipped/non-clipped ratio of fish in counting window at Lower Granite Dam																						
2																							
3	1997																						
4	Non-expanded data from WDFW																						
5																							
6	Lower Granite Dam count of:																						
7	Total			Adult		Jack		Total		Adult		Jack											
8	Chinook			Chinook		Chinook		Chinook		Chinook		Chinook											
9	Spring 3/1 - 6/17			33938		33854		84		23,924		4,360											
10	Summer 6/18 - 8/17			10766		10648		118		5,672		3,357											
11	Fall 8/18 - 12/15			1891		1412		479		517		1,118											
12	Sum of Sp/Su			44704																			
13	Hatchery jacks																						
14	Percentage of fish in category																						
15	Estimated wild spring chinook			2328		2		82		0.85		0.15											
16	Estimated wild summer chinook			3258		57		61		0.63		0.37											
17	total																						
18	5588 60 142																						
19	See sheet "Summary - hatch % identifiable" for percentages used to adjust estimated wild numbers																						
20	1998																						
21	Non-expanded data from WDFW																						
22																							
23	Lower Granite Dam count of:																						
24	Total			Adult		Jack		Total		Adult		Jack											
25	Chinook			Chinook		Chinook		Chinook		Chinook		Chinook											
26	Spring 3/1 - 6/17			9987		9881		106		4,996		3,312											
27	Summer 6/18 - 8/17			4760		4439		321		1,411		2,492											
28	Fall 8/18 - 12/15			3839		1862		1977		2,195		1,089											
29	Sum of Sp/Su			14747		14320		427		Hatchery jacks		Percentage of fish in category											
30	3302 36 70 0.60 0.40 0.60 0.40 0.61 0.39																						
31	Estimated wild summer chinook																						
	2748 113 208 0.36 0.64 0.34 0.66 0.59 0.41																						
	total																						
	6050 149 278																						
	See sheet "Summary - hatch % identifiable" for percentages used to adjust estimated wild numbers																						

	A	B	C
1	name	mass	threat status
2	Rodentia		
3	rat	56	LC
4	mouse	90	LC
5	squirrel	24	CR
6	Primates		
7	gibbon	7000	LC
8	bushbaby	678	EN
9	Chiroptera		
10	myotis	45	NT
11	smaller myotis	48	NT
12	fishing bat	89	VU
13	bulldog bat	33	DD



My workbook has 17 sheets,
but I transposed the data
matrix in sheet 4, at
random, for absolutely no
reason.



The workbook you opened contains automatic links to information in another workbook.

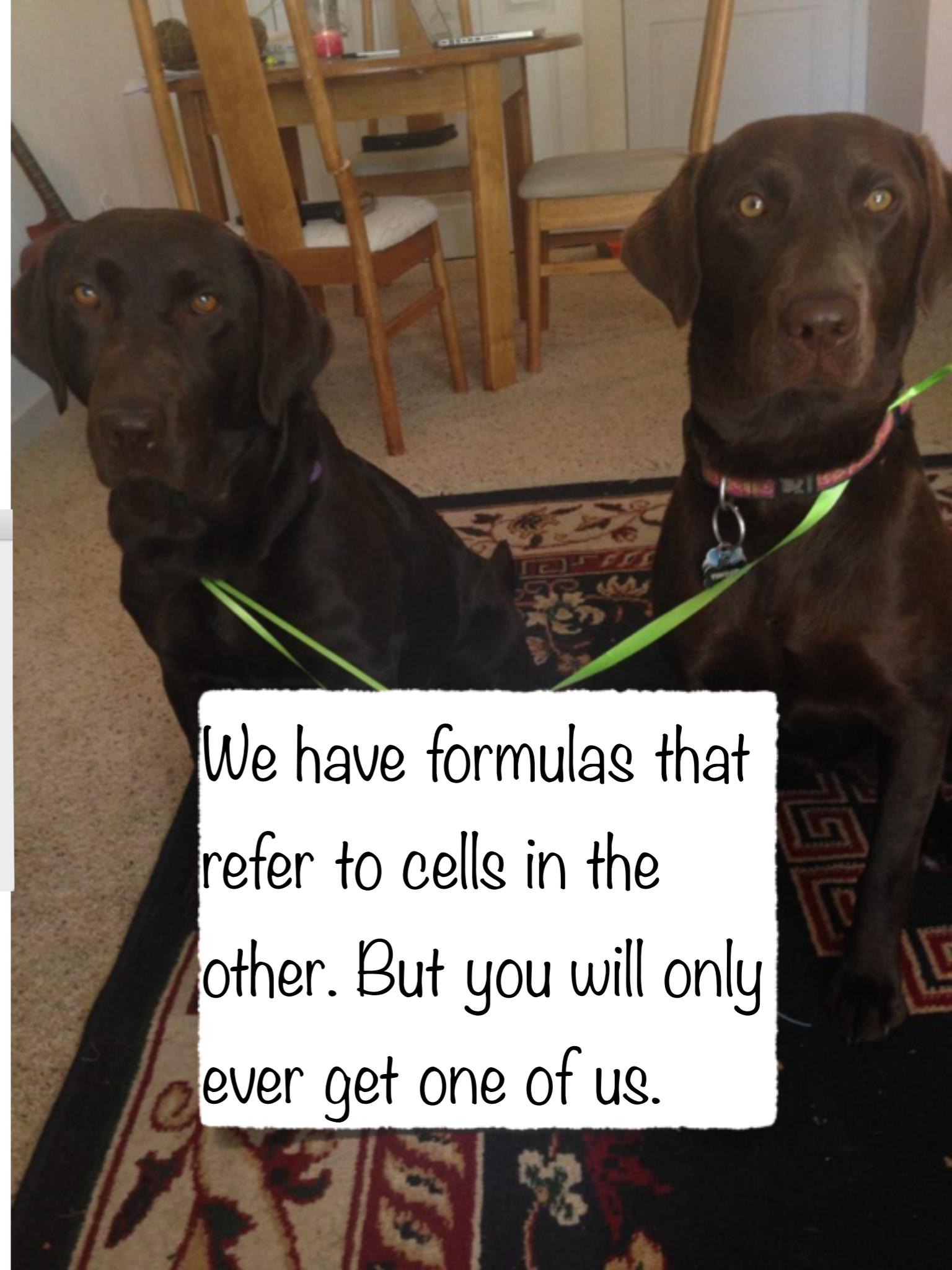
Do you want to update this workbook with changes made to the other workbook?

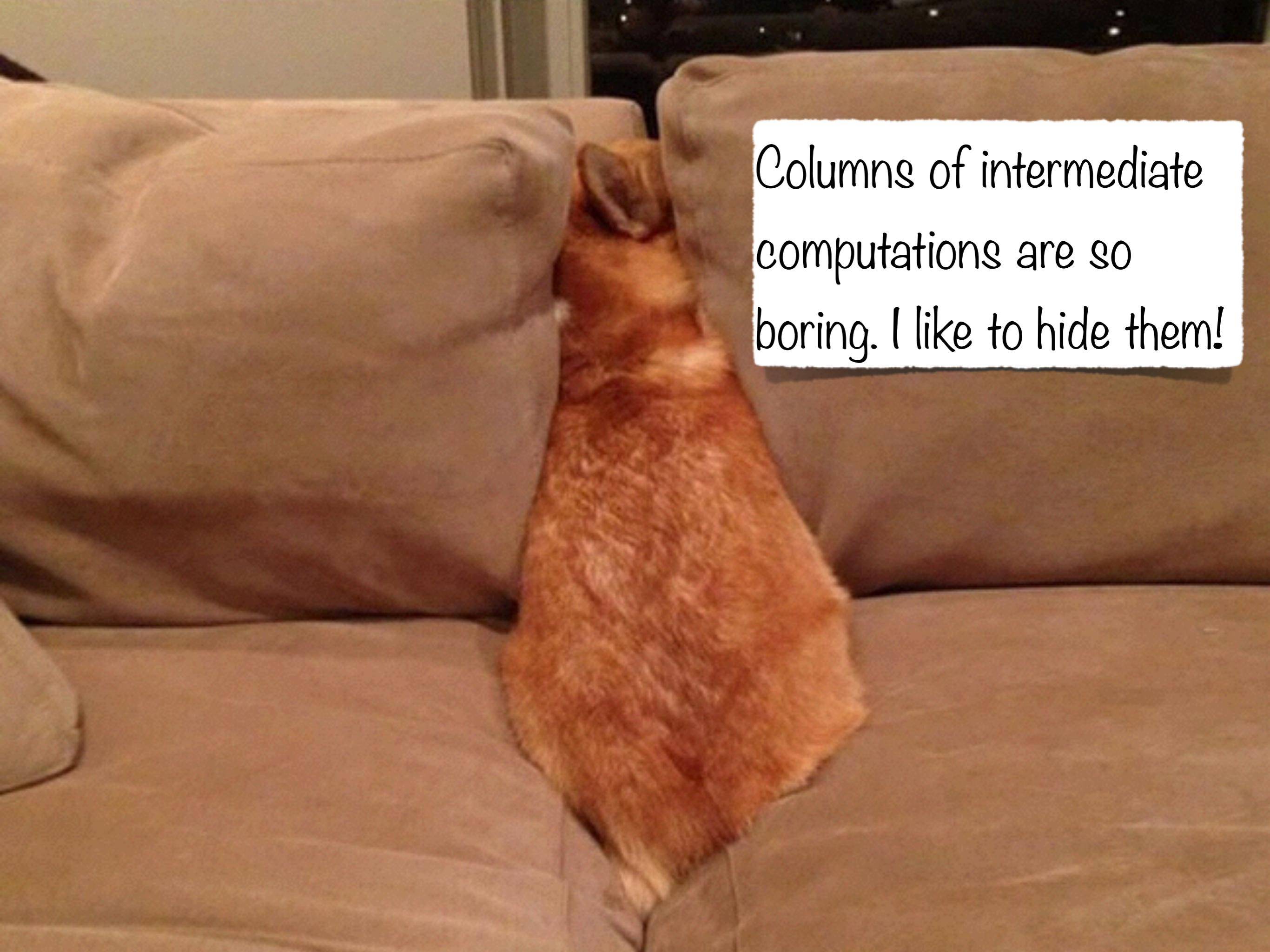
- To update all linked information, click Update. You must have access to all of the linked workbooks.
- To keep the existing information, click Ignore Links.
- To open your workbook and receive more options to which links get updated, click Edit Links.

[Edit Links](#)

[Update](#)

[Ignore Links](#)



A fluffy orange cat with white paws and a white patch on its chest is sitting on a light-colored sofa. It is looking directly at the camera. A white speech bubble with a torn paper effect contains the text.

Columns of intermediate computations are so boring. I like to hide them!

ALGORITHMS BY COMPLEXITY

MORE COMPLEX →

LETPAD QUICKSORT GIT
MERGE SELF-DRIVING
 CAR GOOGLE
 SEARCH BACKEND

SPRAWLING EXCEL SPREADSHEET
BUILT UP OVER 20 YEARS BY A
CHURCH GROUP IN NEBRASKA TO
COORDINATE THEIR SCHEDULING

machine readable

&

human readable

code

can be
machine & human readable

data

can be
machine & human readable



branch: master

lotr / lotr_clean.tsv



file | 684 lines (683 sloc) | 42.64 kb

Open

Edit

Raw

Blame

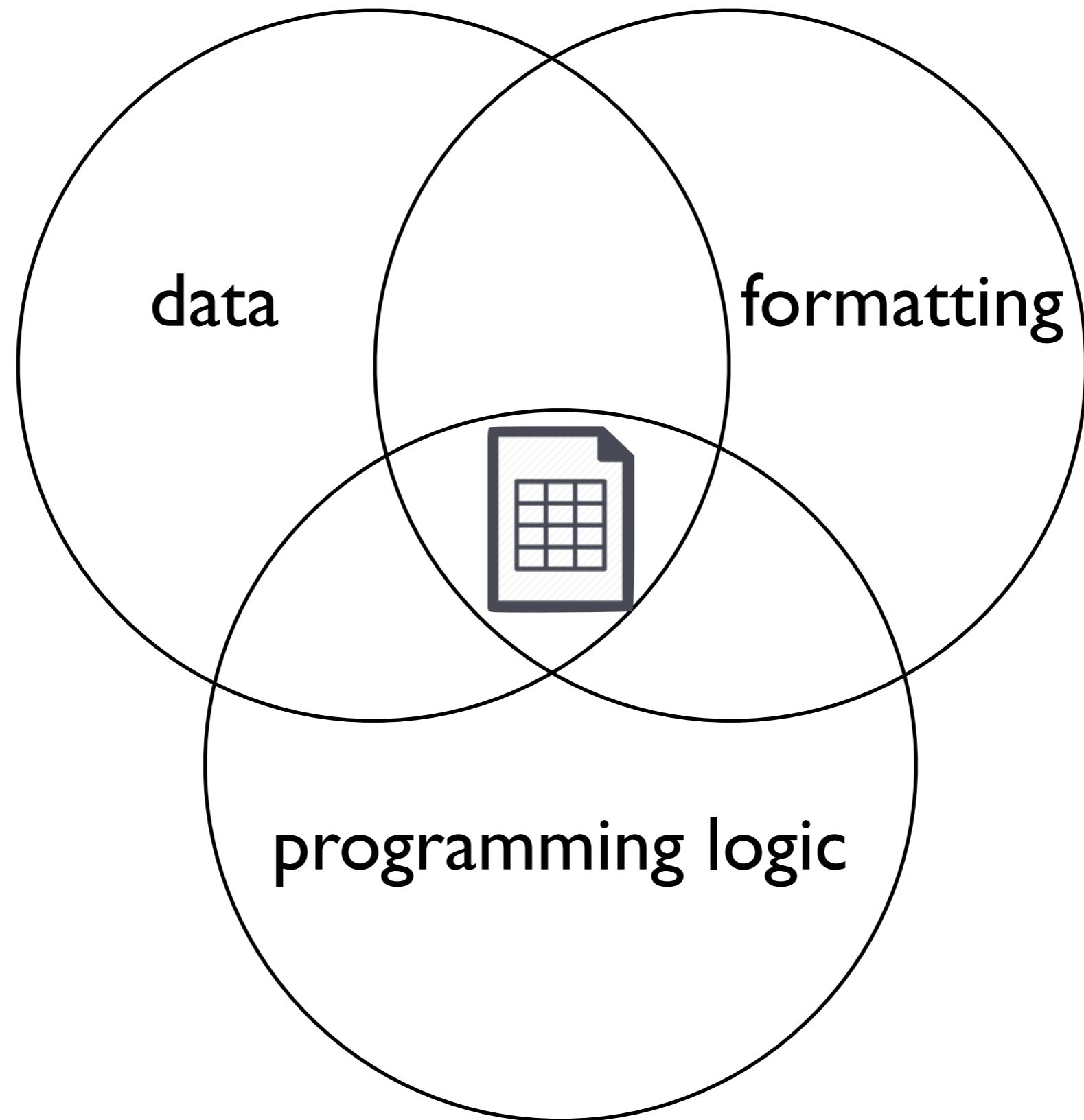
History

Delete

Search this file...

1	Film	Chapter	Character	Race	Words
2	The Fellowship Of The Ring	01: Prologue	Bilbo	Hobbit	4
3	The Fellowship Of The Ring	01: Prologue	Elrond	Elf	5
4	The Fellowship Of The Ring	01: Prologue	Galadriel	Elf	460
5	The Fellowship Of The Ring	02: Concerning Hobbits	Bilbo	Hobbit	214
6	The Fellowship Of The Ring	03: The Shire	Bilbo	Hobbit	70
7	The Fellowship Of The Ring	03: The Shire	Frodo	Hobbit	128
8	The Fellowship Of The Ring	03: The Shire	Gandalf	Wizard	197
9	The Fellowship Of The Ring	03: The Shire	Hobbit Kids	Hobbit	10

a spreadsheet
is often neither
machine nor human readable



what are the problems?

which ones can we solve?

via training

via tooling

be realistic,
be fair,
be precise

let 1,000 flowers bloom!



Two angles on the Spreadsheet Problem:

Create new spreadsheet implementations
that use, e.g., R for computation
and visualization.

Accept spreadsheets as they are.
Create tools to get goodies out and into, e.g., R.
Maybe write back into sheets?



Stencila

Sheets

Spreadsheets are probably the most widely used environment for data analysis and end user programming. Their ubiquity comes from the intuitive simplicity of a live, reactive, see-your-data-and-do-stuff with it interface.

Stencila sheets combine the benefits of the spreadsheet interface with the power of languages like R and Python.

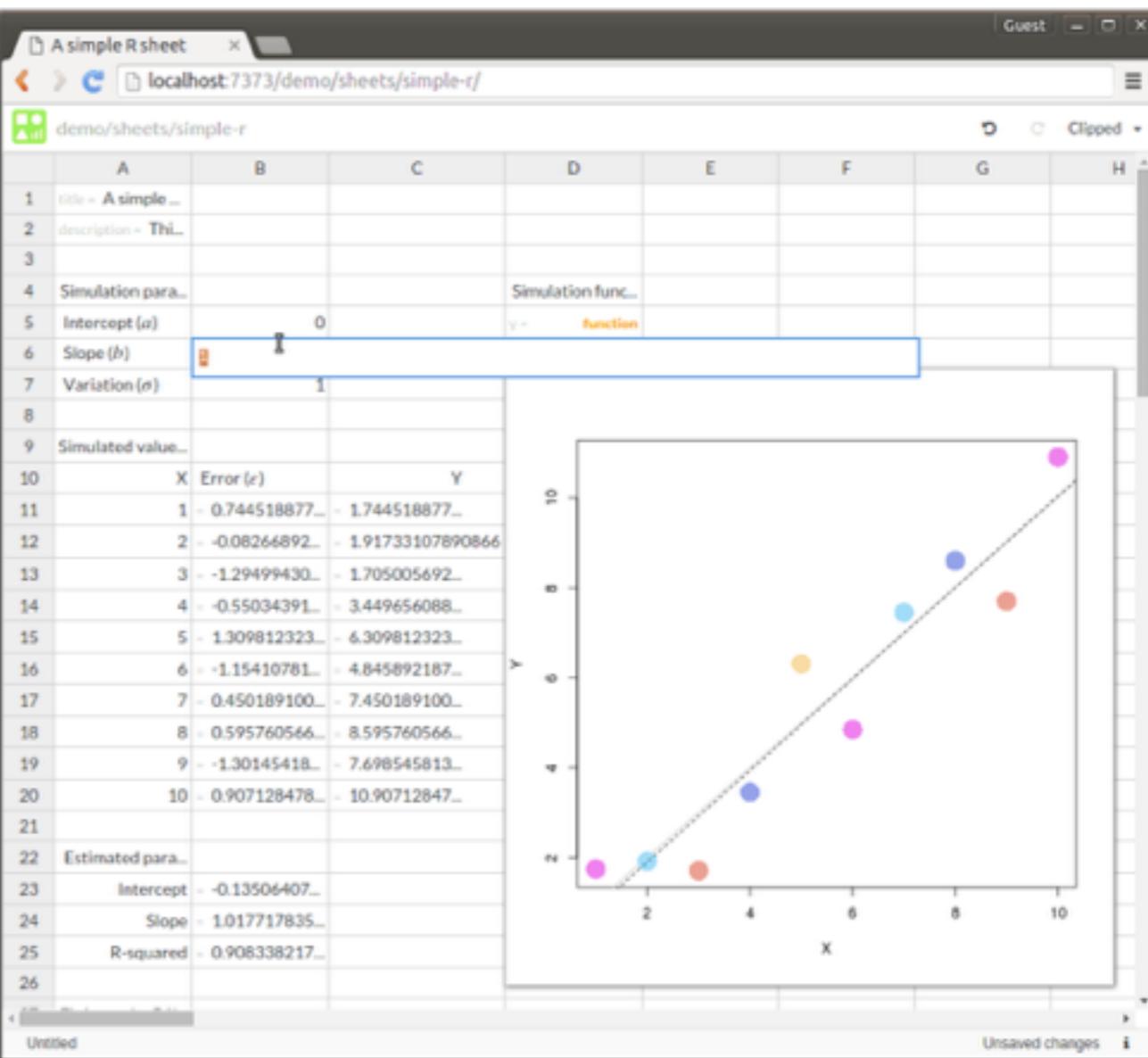
Built from the ground up for transparency, testability and version control, we think Stencila sheets are the next step in the evolution of the spreadsheet.

More:

[Spreadsheets are dead, long live reactive programming environments!](#)

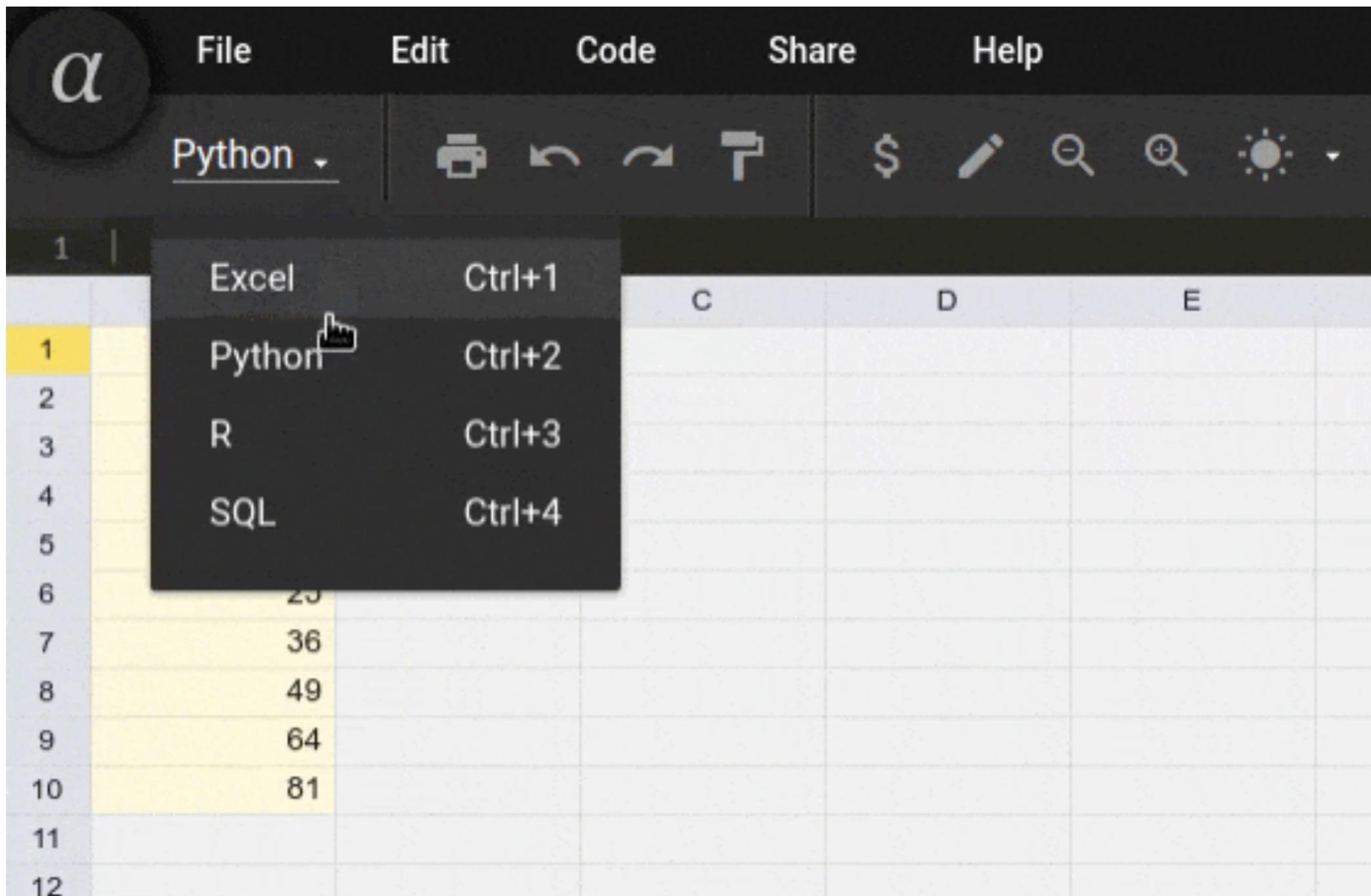
[Getting under Stencila sheets](#)

[A spreadsheet file format for humans](#)



AlphaSheets

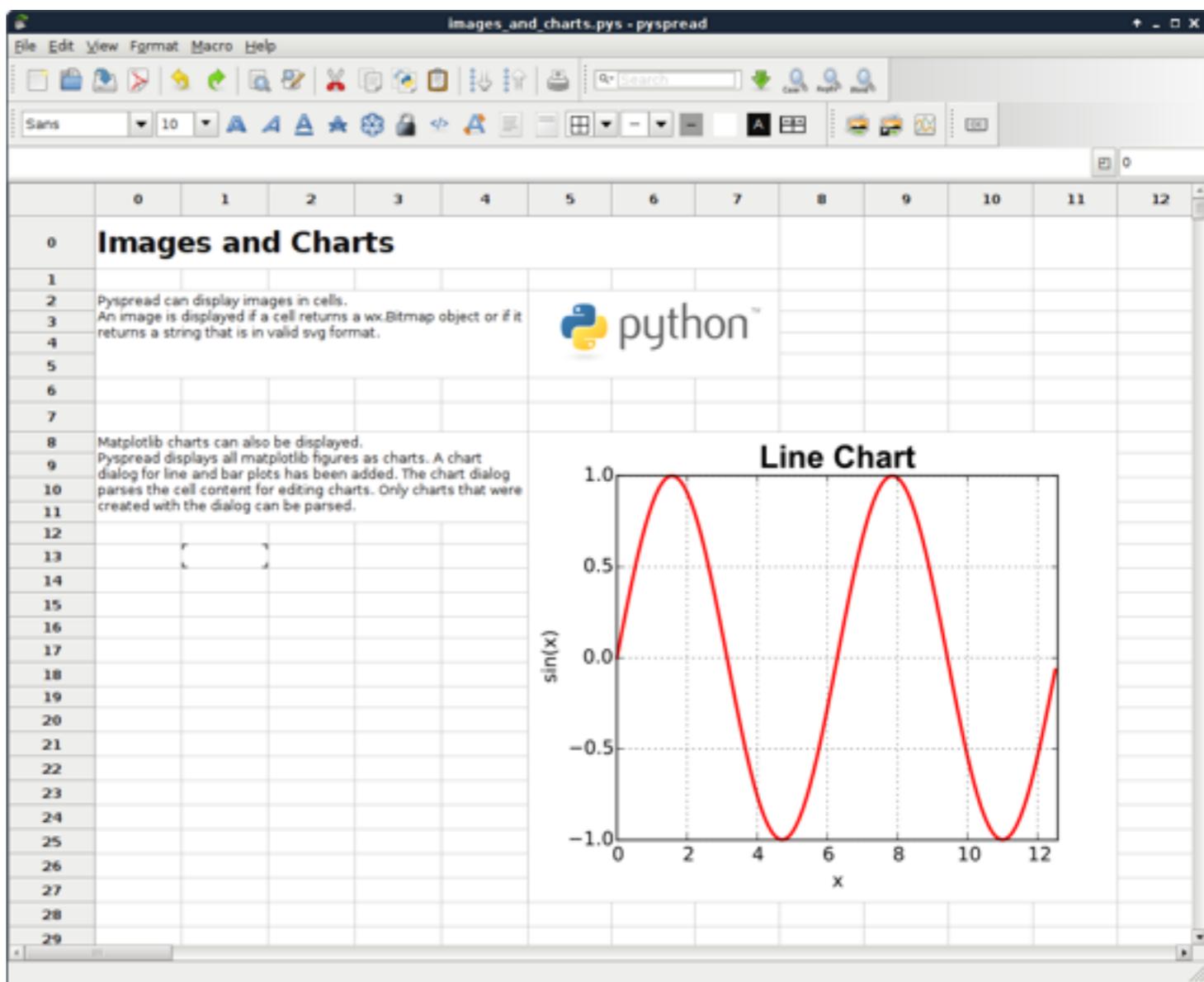
“collaborative, programmable spreadsheets”





pyspread

Python power for your tables

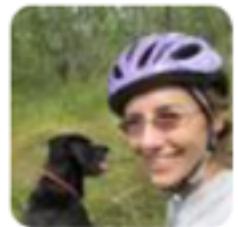


Pyspread is a non-traditional spreadsheet application that is based on and written in the programming language [Python](#).

The goal of pyspread is to be the most pythonic spreadsheet.

Pyspread expects Python expressions in its grid cells, which makes a spreadsheet specific language obsolete. Each cell returns a Python object that can be accessed from other cells. These objects can represent anything including lists or matrices.

Pyspread is free software. It is released under the [GPL v3](#). You can find the source code at [github](#).



Jean Adams

@JeanVAdams

If your collaborator asks, “In what form would you like the data?” you should respond, “In its current form.” via @kwbroman



Enron North America - West Gas

November 9, 2001



ENA - West Gas Contacts

Houston Office

Barry Tycholiz	(713) 853-1587
Kim Ward	(713) 853-0685
Stephanie Miller	(713) 853-1688
Philip Polsky	(713) 853-5181

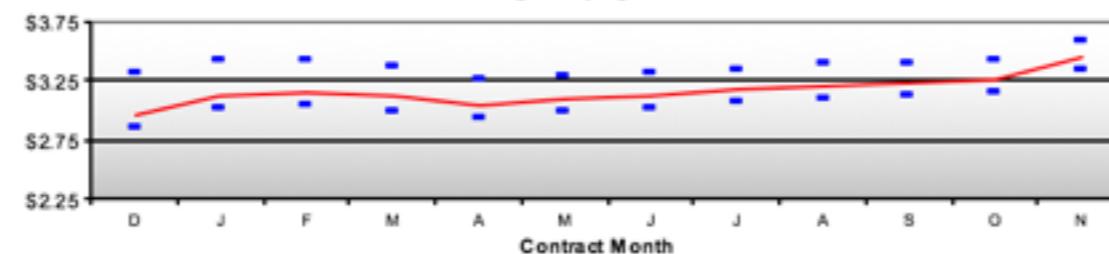
Regional Offices

Mark Whitt	(303) 575-6473	Denver
Paul Lucci	(303) 575-6474	Denver
Tyrell Harrison	(303) 575-6478	Denver
Dave Fuller	(503) 464-3732	Portland

Forward Prices (US\$/MMBtu)

NYMEX	
	SETTLE
Cash	
ROM	
Dec-01	2.960 0.090
Dec-01 to Mar-02	3.088 0.083
Apr-02 to Oct-02	3.166 0.084
Nov-02 to Mar-03	3.651 0.090
One Year Strip*	3.165 0.084

Forward NYMEX Strip
with trailing 10-day highs/lows



IF NWPL Rocky Mountains

Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.890	1.910		
2.060	2.080		
2.395	2.415	(0.565)	(0.545)
2.594	2.614	(0.494)	(0.474)
2.581	2.601	(0.585)	(0.565)
3.356	3.376	(0.295)	(0.275)
2.634	2.654	(0.530)	(0.510)

IF CIG Rocky Mountains

Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.940	1.960		
1.960	1.980		
2.345	2.365	(0.615)	(0.595)
2.548	2.568	(0.540)	(0.520)
2.471	2.491	(0.695)	(0.675)
3.311	3.331	(0.340)	(0.320)
2.551	2.571	(0.614)	(0.594)

IF EL Paso Permian

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.375	2.395		
2.420	2.440		
2.700	2.720	(0.260)	(0.240)
2.855	2.875	(0.233)	(0.213)
3.009	3.029	(0.158)	(0.138)
3.499	3.519	(0.153)	(0.133)
2.982	3.002	(0.182)	(0.162)

IF EL Paso San Juan

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.450	2.470		
2.350	2.370		
2.560	2.580	(0.400)	(0.380)
2.743	2.763	(0.345)	(0.325)
2.801	2.821	(0.365)	(0.345)
3.421	3.441	(0.230)	(0.210)
2.817	2.837	(0.347)	(0.327)

AECO / NIT

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.376	2.396		
2.398	2.418		
2.552	2.572	(0.408)	(0.388)
2.616	2.636	(0.472)	(0.452)
2.661	2.681	(0.505)	(0.485)
3.216	3.236	(0.435)	(0.415)
2.676	2.696	(0.488)	(0.468)

IF NWPL Canadian Border (Sumas)

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.480	2.500		
2.460	2.480		
2.800	2.820	(0.160)	(0.140)
2.892	2.912	(0.196)	(0.176)
2.796	2.816	(0.370)	(0.350)
3.706	3.726	0.055	0.075
2.880	2.900	(0.285)	(0.265)

IF PEPL TX-OK

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.530	2.550		
2.530	2.550		
2.828	2.848	(0.133)	(0.113)
2.958	2.978	(0.130)	(0.110)
3.046	3.066	(0.120)	(0.100)
3.531	3.551	(0.120)	(0.100)
3.041	3.061	(0.123)	(0.103)

Tidying an untidyable dataset

David Robinson

April 7, 2016

Jenny Bryan
@JennyBryan
@JennyBryan check out this beauty pic.twitter.com/U5EVm11mj

6 Apr

Simply Statistics
@simplystats

[Follow](#)

@JennyBryan I want to see @drob figure out how to turn this into a tidy data set. I'd pay money to watch it in realtime....

1:44 PM - 6 Apr 2016

2 16

Challenge accepted. Here's the data Jenny sent, which can be downloaded [here](#):

~150 lines of code later . . .

1.89	1.91	-0.5	-0.5
2.06	2.08	-0.4	-0.4
2.38	2.41	-0.5	-0.5
2.59	2.61	-0.5	-0.5
2.58	2.60	-0.5	-0.5
3.35	3.37	-0.5	-0.5
3.35	3.37	-0.5	-0.5
2.63	2.65	-0.5	-0.5

1.94	1.96	-0.6	-0.5
1.96	1.98	-0.6	-0.5
2.34	2.36	-0.6	-0.5
2.54	2.56	-0.5	-0.5
2.47	2.49	-0.6	-0.6
3.31	3.33	-0.3	-0.3
2.55	2.57	-0.6	-0.5

2.37	2.39	-0.2	-0.2
2.42	2.44	-0.2	-0.2
2.7	2.72	-0.2	-0.2
2.85	2.87	-0.1	-0.1
3.00	3.02	-0.1	-0.1
3.49	3.51	-0.1	-0.1
2.98	3.00	-0.1	-0.1

2.45	2.47	-0.4	-0.3
2.35	2.37	-0.3	-0.3
2.56	2.58	-0.3	-0.3
2.74	2.76	-0.3	-0.3
2.80	2.82	-0.3	-0.3
3.42	3.44	-0.2	-0.2
2.81	2.83	-0.3	-0.3

BID
OFFER

2.37	2.39	-0.4	-0.3
2.39	2.41	-0.4	-0.3
2.55	2.57	-0.4	-0.3
2.61	2.63	-0.4	-0.4
2.66	2.68	-0.5	-0.4
3.21	3.23	-0.4	-0.4
2.67	2.69	-0.4	-0.4

2.48	2.5	-0.1	-0.1
2.46	2.48	-0.1	-0.1
2.8	2.82	-0.1	-0.1
2.89	2.91	-0.3	-0.3
2.79	2.81	-0.3	-0.3
3.70	3.72	5.5E	7.49
2.87	2.89	-0.2	-0.2

2.52	2.54	-0.1	-0.1
2.52	2.54	-0.1	-0.1
2.82	2.84	-0.1	-0.1
2.95	2.97	-0.1	-0.1
3.04	3.06	-0.1	-0.1
3.53	3.55	-0.1	-0.1
3.04	3.06	-0.1	-0.1

2.58	2.6	-0.1	-0.1
2.5	2.52	-0.1	-0.1
2.79	2.81	-0.1	-0.1
2.94	2.96	-0.1	-0.1
3.22	3.24	5.85	7.85
3.74	3.76	0.09	0.11
3.15	3.17	-5.0	1.49

2.54	2.57	-0.1	-0.1
2.48	2.5	-0.1	-0.1
2.78	2.80	-0.1	-0.1
2.91	2.93	-0.1	-0.1
3.04	3.06	-0.1	-0.1
3.72	3.74	7.00	0.09
3.03	3.05	-0.1	-0.1

2.57	2.59	-0.0	-0.0
2.52	2.54	-0.2	-0.2
2.88	2.90	-0.2	-0.2
3.01	3.03	-0.2	-0.2
3.26	3.28	9.50	0.11
3.95	3.97	0.30	0.32
3.21	3.23	5.12	7.12



readxl: [CRAN](#), [GitHub](#)

openxlsx: [CRAN](#), [GitHub](#)

XLConnect: [CRAN](#), [GitHub](#)

xlsx: [CRAN](#), [GitHub](#)

gdata: [CRAN](#), [R-Forge](#)

... and more

What do we want?

no tricky dependency ... no Java
agnostic re: Excel, Google Sheet, ill-formed csv

expose

(unformatted) data

(unevaluated) formulas

formatting

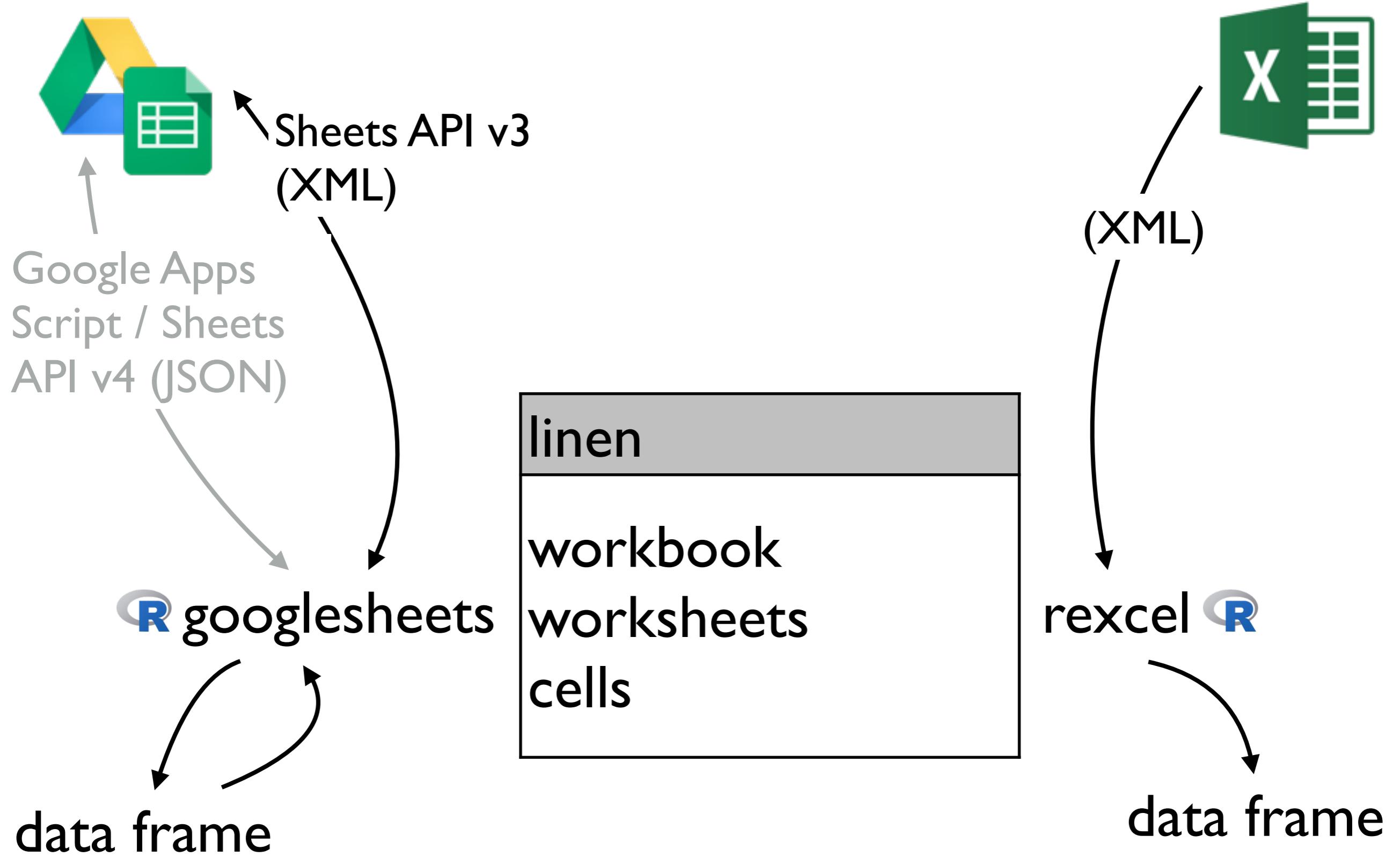
detect / propose views

handle merged cells, weird headers

How are we doing it?

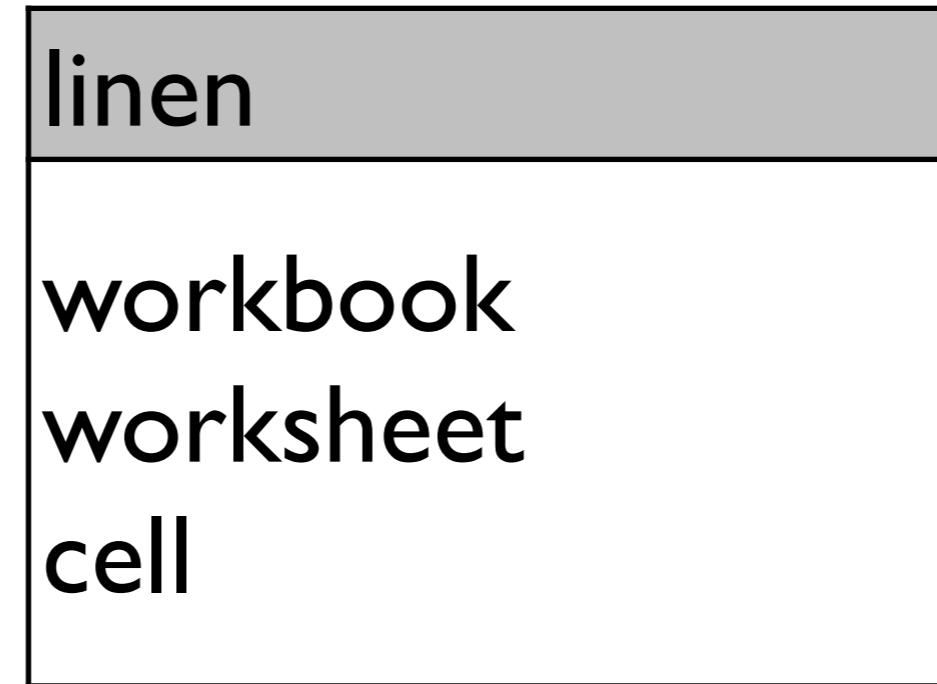
define the `linen` object = spreadsheet receptacle
document meta-data
worksheet meta-data
cell data, broadly defined

`rexcel` & `googlesheets` create `linen` objects
simple? return a data frame!
not? expose `linen` object for more processing ...





R googlesheets



Rexcel R

jailbreakr

multiple views, data frames
unformatted data, formatting
unevaluated formulas
figures?



A

B C D E F G H I J K L M N O P Q R S T U

Enron North America - West Gas

November 9, 2001



ENA - West Gas Contacts

Houston Office

Barry Tycholiz	(713) 853-1587
Kim Ward	(713) 853-0685
Stephanie Miller	(713) 853-1688
Philip Polsky	(713) 853-5181

Regional Offices

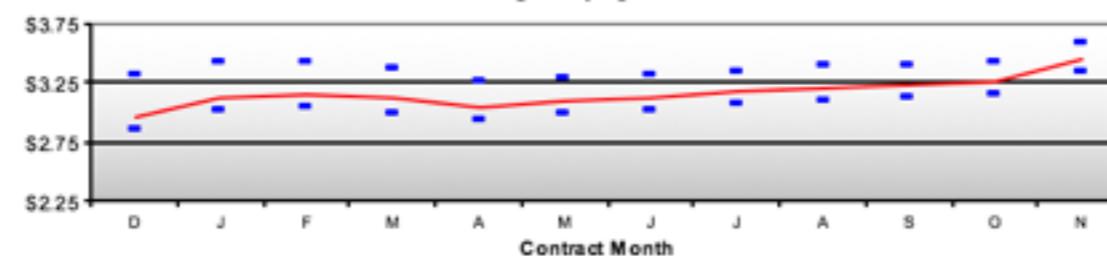
Mark Whitt	(303) 575-6473	Denver
Paul Lucci	(303) 575-6474	Denver
Tyrell Harrison	(303) 575-6478	Denver
Dave Fuller	(503) 464-3732	Portland

Forward Prices (US\$/MMBtu)

NYMEX

	SETTLE	Δ
Cash		
ROM		
Dec-01	2.960	0.090
Dec-01 to Mar-02	3.088	0.083
Apr-02 to Oct-02	3.166	0.084
Nov-02 to Mar-03	3.651	0.090
One Year Strip*	3.165	0.084

Forward NYMEX Strip
with trailing 10-day highs/lows



IF NWPL Rocky Mountains

Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.890	1.910		
2.060	2.080		
2.395	2.415	(0.565)	(0.545)
2.594	2.614	(0.494)	(0.474)
2.581	2.601	(0.585)	(0.565)
3.356	3.376	(0.295)	(0.275)
2.634	2.654	(0.530)	(0.510)

IF CIG Rocky Mountains

Fixed Price		Basis		
	BID	OFFER	BID	OFFER
Cash	1.940	1.960		
ROM	1.960	1.980		
Dec-01	2.345	2.365	(0.615)	(0.595)
Dec-01 to Mar-02	2.548	2.568	(0.540)	(0.520)
Apr-02 to Oct-02	2.471	2.491	(0.695)	(0.675)
Nov-02 to Mar-03	3.311	3.331	(0.340)	(0.320)
One Year Strip*	2.551	2.571	(0.614)	(0.594)

IF EL Paso Permian

Fixed Price		Basis		
	BID	OFFER	BID	OFFER
	2.375	2.395		
	2.420	2.440		
	2.700	2.720	(0.260)	(0.240)
	2.855	2.875	(0.233)	(0.213)
	3.009	3.029	(0.158)	(0.138)
	3.499	3.519	(0.153)	(0.133)
	2.982	3.002	(0.182)	(0.162)

IF EL Paso San Juan

Fixed Price		Basis		
	BID	OFFER	BID	OFFER
	2.450	2.470		
	2.350	2.370		
	2.560	2.580	(0.400)	(0.380)
	2.743	2.763	(0.345)	(0.325)
	2.801	2.821	(0.365)	(0.345)
	3.421	3.441	(0.230)	(0.210)
	2.817	2.837	(0.347)	(0.327)

AECO / NIT

Fixed Price		Basis		
	BID	OFFER	BID	OFFER
Cash	2.376	2.396		
ROM	2.398	2.418		
Dec-01	2.552	2.572	(0.408)	(0.388)
Dec-01 to Mar-02	2.616	2.636	(0.472)	(0.452)
Apr-02 to Oct-02	2.661	2.681	(0.505)	(0.485)
Nov-02 to Mar-03	3.216	3.236	(0.435)	(0.415)
One Year Strip*	2.676	2.696	(0.488)	(0.468)

IF NWPL Canadian Border (Sumas)

Fixed Price		Basis		
	BID	OFFER	BID	OFFER
	2.480	2.500		
	2.460	2.480		
	2.800	2.820	(0.160)	(0.140)
	2.892	2.912	(0.196)	(0.176)
	2.796	2.816	(0.370)	(0.350)
	3.706	3.726	0.055	0.075
	2.880	2.900	(0.285)	(0.265)

IF PEPL TX-OK

Fixed Price		Basis		
	BID	OFFER	BID	OFFER
	2.530	2.550		
	2.530	2.550		
	2.828	2.848	(0.133)	(0.113)
	2.958	2.978	(0.130)	(0.110)
	3.046	3.066	(0.120)	(0.100)
	3.531	3.551	(0.120)	(0.100)
	3.041	3.061	(0.123)	(0.103)



Enron North America - West Gas

November 9, 2001



ENA - West Gas Contacts

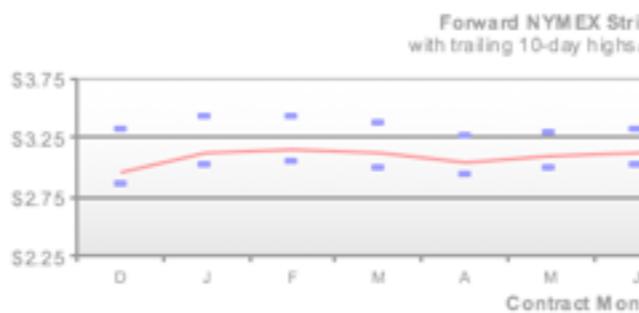
Houston Office

Barry Tycholiz	(713) 853-1587
Kim Ward	(713) 853-0685
Stephanie Miller	(713) 853-1688
Philip Polsky	(713) 853-5181

Regional Offices

Mark Whitt	(303) 575-6473	Denver
Paul Lucci	(303) 575-6474	Denver
Tyrell Harrison	(303) 575-6478	Denver
Dave Fuller	(503) 464-3732	Portland

Forward Prices (US\$/MMBtu)



SETTLE	Δ
2.960	0.090
3.088	0.083
3.166	0.084
3.651	0.090
3.165	0.084

IF CIG Rocky Mountains			
Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.940	1.960		
1.960	1.980		
2.345	2.365	(0.615)	(0.595)
2.548	2.568	(0.540)	(0.520)
2.471	2.491	(0.695)	(0.675)
3.311	3.331	(0.340)	(0.320)
2.551	2.571	(0.614)	(0.594)

AECO / NIT			
Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.376	2.396		
2.398	2.418		
2.552	2.572	(0.408)	(0.388)
2.616	2.636	(0.472)	(0.452)
2.661	2.681	(0.505)	(0.485)
3.216	3.236	(0.435)	(0.415)
2.676	2.696	(0.488)	(0.468)

AC

AD

AE

AF

AG

115

116 **Congratulations to Dan Foulston
117 of Canadian Western Natural Gas!!!**

118 **He is the winner of our weekend getaway
119 to Lake Louise. Dan's guess of \$1.375 C/GJ
120 was the closest to the average Aeco
121 day gas price for the first
122 15 days of December.**

113:

114:

115:

116:

117:

118:

119:

120:

121:

122:

d

d

d

d

d

d

d



Enron North America - West Gas

November 9, 2001



ENA - West Gas Contacts

Houston Office

Barry Tycholiz	(713) 853-1587
Kim Ward	(713) 853-0685
Stephanie Miller	(713) 853-1688
Philip Polsky	(713) 853-5181

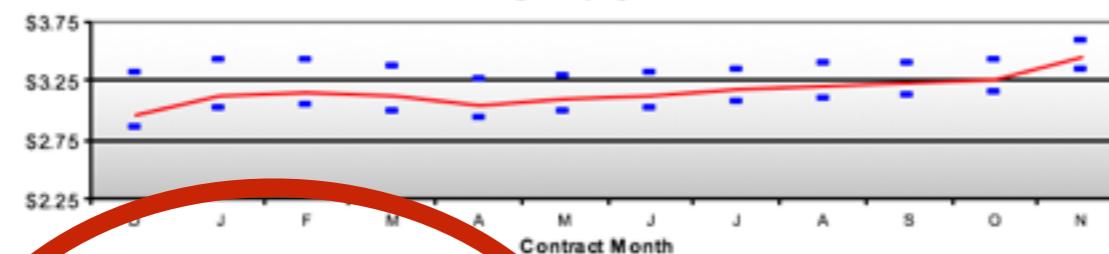
Regional Offices

Mark Whitt	(303) 575-6473	Denver
Paul Lucci	(303) 575-6474	Denver
Tyrell Harrison	(303) 575-6478	Denver
Dave Fuller	(503) 464-3732	Portland

Forward Prices (US\$/MMBtu)

NYMEX	
	SETTLE
Cash	
ROM	
Dec-01	2.960 0.090
Dec-01 to Mar-02	3.088 0.083
Apr-02 to Oct-02	3.166 0.084
Nov-02 to Mar-03	3.651 0.090
One Year Strip*	3.165 0.084

Forward NYMEX Strip
with trailing 10-day highs/lows



IF NWPL Rocky Mountains

Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.890	1.910		
2.060	2.080		
2.395	2.415	(0.565)	(0.545)
2.594	2.614	(0.494)	(0.474)
2.581	2.601	(0.585)	(0.565)
3.356	3.376	(0.295)	(0.275)
2.634	2.654	(0.530)	(0.510)

IF CIG Rocky Mountains

Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.940	1.960		
1.960	1.980		
2.345	2.365	(0.615)	(0.595)
2.548	2.568	(0.540)	(0.520)
2.471	2.491	(0.695)	(0.675)
3.311	3.331	(0.340)	(0.320)
2.551	2.571	(0.614)	(0.594)

IF EL Paso Permian

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.375	2.395		
2.420	2.440		
2.700	2.720	(0.260)	(0.240)
2.855	2.875	(0.233)	(0.213)
3.009	3.029	(0.158)	(0.138)
3.499	3.519	(0.153)	(0.133)
2.982	3.002	(0.182)	(0.162)

IF EL Paso San Juan

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.450	2.470		
2.350	2.370		
2.560	2.580	(0.400)	(0.380)
2.743	2.763	(0.345)	(0.325)
2.801	2.821	(0.365)	(0.345)
3.421	3.441	(0.230)	(0.210)
2.817	2.837	(0.347)	(0.327)

AECO / NIT

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.376	2.396		
2.398	2.418		
2.552	2.572	(0.408)	(0.388)
2.616	2.636	(0.472)	(0.452)
2.661	2.681	(0.505)	(0.485)
3.216	3.236	(0.435)	(0.415)
2.676	2.696	(0.488)	(0.468)

IF NWPL Canadian Border (Sumas)

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.480	2.500		
2.460	2.480		
2.800	2.820	(0.160)	(0.140)
2.892	2.912	(0.196)	(0.176)
2.796	2.816	(0.370)	(0.350)
3.706	3.726	0.055	0.075
2.880	2.900	(0.285)	(0.265)

IF PEPL TX-OK

Fixed Price		Basis	
BID	OFFER	BID	OFFER
2.530	2.550		
2.530	2.550		
2.828	2.848	(0.133)	(0.113)
2.958	2.978	(0.130)	(0.110)
3.046	3.066	(0.120)	(0.100)
3.531	3.551	(0.120)	(0.100)
3.041	3.061	(0.123)	(0.103)



Enron North America - West Gas

November 9, 2001



ENA - West Gas Contacts

Houston Office

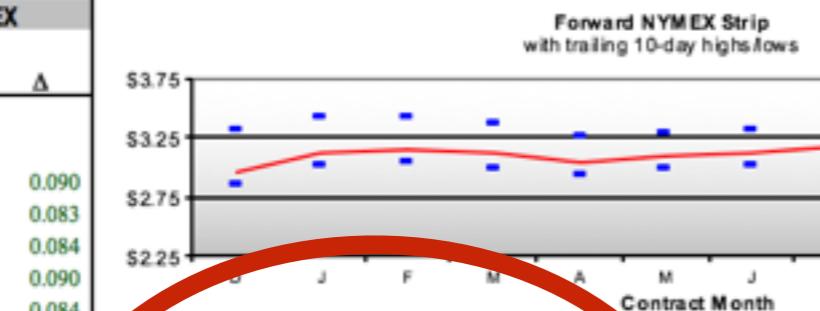
Barry Tycholiz	(713) 853-1587
Kim Ward	(713) 853-0685
Stephanie Miller	(713) 853-1688
Philip Polsky	(713) 853-5181

Regional Offices

Mark Whitt	(303) 575-6473	Denver
Paul Lucci	(303) 575-6474	Denver
Tyrell Harrison	(303) 575-6478	Denver
Dave Fuller	(503) 464-3732	Portland

Forward Prices (US\$/MMBtu)

SETTLE	Δ
2.960	0.090
3.088	0.083
3.166	0.084
3.651	0.090
3.165	0.084



IF NWPL Rocky Mountains

```
> views <- jailbreakr::split_sheet(sheet)
> length(views)
```

[1] 40

```
[1] > head(views, 3)
```

[1]

<worksheet_view: 2 x 20 (of 235 x 99)>

: BCDEFGHIJKLMNOPQRSTUVWXYZ

1: a←←←←←←←←←←
2: Ø←←←←←←←←←←

IF CIG Rocky Mountains			
Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.940	1.960		
1.960	1.980		
2.345	2.365	(0.615)	(0.595)
2.548	2.568	(0.540)	(0.520)
2.471	2.491	(0.695)	(0.675)
3.311	3.331	(0.340)	(0.320)
2.551	2.571	(0.614)	(0.594)

Fixed
BID
2.375
2.420
2.700
2.855
3.009
3.499
2.982

Fixed Price Basis			
BID	OFFER	BID	OFFER
2.376	2.396		
2.398	2.418		
2.552	2.572	(0.408)	(0.388)
2.616	2.636	(0.472)	(0.452)
2.661	2.681	(0.505)	(0.485)
3.216	3.236	(0.435)	(0.415)
2.676	2.696	(0.488)	(0.468)

IF NW
Fixed P
BID
2.480
2.460
2.800
2.892
2.796
3.706
2.880

[[2]]

worksheet_view: 1 x 20

[3]

<worksheet view: 1 x 20 (of 235 x 99)>

: BCDEFGHIJKLMNOPQRSTUVWXYZ

IF CIG Rocky Mountains

Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.940	1.960		
1.960	1.980		
2.345	2.365	(0.615)	(0.595)
2.548	2.568	(0.540)	(0.520)
2.471	2.491	(0.695)	(0.675)
3.311	3.331	(0.340)	(0.320)
2.551	2.571	(0.614)	(0.594)

```
> x2 <- jailbreakr::split_metadata(x1, FALSE, min_data_block=1)
> x3 <- jailbreakr::split_headers_apply(x2, 2)
> x3
<worksheet_view: 7 x 4 (of 235 x 99)>
  (with headers)
  (with data: metadata)
    : GHIJ
  28: 00
  29: 00
  30: 0000
  31: 0000
  32: 0000
  33: 0000
  34: 0000
```

IF CIG Rocky Mountains

Fixed Price		Basis	
BID	OFFER	BID	OFFER
1.940	1.960		
1.960	1.980		
2.345	2.365	(0.615)	(0.595)
2.548	2.568	(0.540)	(0.520)
2.471	2.491	(0.695)	(0.675)
3.311	3.331	(0.340)	(0.320)
2.551	2.571	(0.614)	(0.594)

```
> x3$values()
   Fixed Price:BID Fixed Price:OFFER Basis:BID Basis:OFFER
[1,] 1.94           1.96          NA          NA
[2,] 1.96           1.98          NA          NA
[3,] 2.345          2.365         -0.615      -0.595
[4,] 2.54775        2.56775       -0.54       -0.52
[5,] 2.471286       2.491286      -0.695      -0.675
[6,] 3.311          3.331         -0.34       -0.32
[7,] 2.55075        2.57075       -0.61375    -0.59375
```



rsheets

<https://github.com/rsheets>

Repositories People 2 Teams 0 Settings

Filters ▾

Find a repository...

New repository

rexcel

Extracts spreadsheet data from Excel workbooks and puts into linen format

Updated 8 days ago

R ★ 6 ⚡ 1

enron_corpus_linen

Code to process the enron corpus

Updated 9 days ago

R ★ 0 ⚡ 0

cellranger

Helper functions to work with spreadsheets and the "A1:D10" style of cell range specification

Updated on May 25

R ★ 12 ⚡ 3

linen

General representation of spreadsheet data, plus some limited low-level operations on that data

Updated on May 13

R ★ 5 ⚡ 0

README

Updated on May 3

★ 0 ⚡ 0

jailbreakr

Get out of Excel free.

Updated on May 3

R ★ 3 ⚡ 0

People

2 >



jennybc

Jennifer (Jenny) Bryan



richfitz

Rich FitzJohn

Invite someone



R `googlesheets`

raw
object data
frame

R `linen`
`cellranger`



R `jailbreakr`

multiple
data frames

formulas,
formatting,
figures?



`rexcel` **R**

raw
object data
frame



bonus content

googlesheets





gs-test-formula-formatting



File Edit View Insert Format Data Tools Add-ons Help



View only

	A	B	C	D	E	F
1	integer	number_formatted	number_rounded	character	formula	formula_formatted
2	123456	654,321		1.23	one	Google
3	345678	12.34%		2.35		52.63%
4	234567	1.23E+09		3.46	three	Google
5		3 1/7		4.57	four	\$A\$1
6	567890	\$0.36		5.68	five	317,898

	A	B	C	D	E	F
1	integer	number_formatted	number_rounded	character	formula	formula_formatted
2	123456	654,321	1.23	one	Google	3.18E+05
3	345678	12.34%	2.35		1,271,591.00	52.63%
4	234567	1.23E+09	3.46	three	Google	0.22
5		3 1/7	4.57	four	\$A\$1	123,456.00
6	567890	\$0.36	5.68	five		317,898

```
gs_ff() %>%
  gs_read() %>%
  select(-integer)
```

#> Accessing worksheet titled 'Sheet1'.
#> No encoding supplied: defaulting to UTF-8.
#> Source: local data frame [5 x 5]

#>

	number_formatted	number_rounded	character	formula	formula_formatted
#>	(chr)	(dbl)	(chr)	(chr)	(chr)
#> 1	654,321	1.23	one	Google	3.18E+05
#> 2	12.34%	2.35	NA	1,271,591.00	52.63%
#> 3	1.23E+09	3.46	three	NA	0.22
#> 4	3 1/7	4.57	four	\$A\$1	123,456.00
#> 5	\$0.36	5.68	five	NA	317,898

default read does not necessarily
give you what you want with
numeric formatting and formulas

	A	B	C	D	E	F
1	integer	number_formatted	number_rounded	character	formula	formula_formatted
2	123456	654,321	1.23	one	Google	3.18E+05
3	345678	12.34%	2.35		1,271,591.00	52.63%
4	234567	1.23E+09	3.46	three	Google	0.22
5		3 1/7	4.57	four	\$A\$1	123,456.00
6	567890	\$0.36	5.68	five		317,898

```

cf <- gs_read_cellfeed(gs_ff())
cf %>%
  filter(row > 1, col == 2) %>%
  select(value, input_value, numeric_value) %>%
  readr::type_convert()
#> #> <tibble [5 x 3]>
#> #>   value  input_value numeric_value
#> #>   <chr>    <dbl>      <dbl>
#> 1 654,321 6.543210e+05 6.543210e+05
#> 2 12.34% 1.234000e+01 1.234000e-01
#> 3 1.23E+09 1.234568e+09 1.234568e+09
#> 4 3 1/7 3.141593e+00 3.141593e+00
#> 5 $0.36 3.600000e-01 3.600000e-01

```

gs-test-formula-formatting

File Edit View Insert Format Data Tools Add-ons Help

View only

fx | integer

	A	B	C	D	E	F
1	integer	number_formatted	number_rounded	character	formula	formula_formatted
2	123456	654,321	1.23	one	Google	3.18E+05
3	345678	12.34%	2.35		1,271,591.00	52.63%
4	234567	1.23E+09	3.46	three	Google	0.22
5		3 1/7	4.57	four	\$A\$1	123,456.00
6	567890	\$0.36	5.68	five		317,898

```
cf <- gs_read_cellfeed(gs_ff())
cf %>%
  filter(row > 1, col == 3) %>%
  select(value, input_value, numeric_value) %>%
  readr::type_convert()
#> #> <tibble [5 x 3]>
#>   value input_value numeric_value
#>   <dbl>      <dbl>      <dbl>
#> 1  1.23      1.2345    1.2345
#> 2  2.35      2.3456    2.3456
#> 3  3.46      3.4567    3.4567
#> 4  4.57      4.5678    4.5678
#> 5  5.68      5.6789    5.6789
```

	A	B	C	D	E	F
1	integer	number_formatted	number_rounded	character	formula	formula_formatted
2	123456	654,321	1.23	one	Google	3.18E+05
3	345678	12.34%	2.35		1,271,591.00	52.63%
4	234567	1.23E+09	3.46	three	Google	0.22
5		3 1/7	4.57	four	\$A\$1	123,456.00
6	567890	\$0.36	5.68	five		317,898

```

cf <- gs_read_cellfeed(gs_ff())
cf %>%
  filter(row > 1, col == 5) %>%
  select(value, input_value, numeric_value) %>%
  mutate(input_value = substr(input_value, 1, 43)) %>%
  readr::type_convert()
#> #<tibble [5 x 3]>
#> #>       value
#> #>       <chr>           input_value numeric_value
#> #>       <chr>           <chr>          <dbl>
#> 1     Google =HYPERLINK("http://www.google.com/","Google") NA
#> 2   1,271,591.00 =sum(R[-1]C[-4]:R[3]C[-4]) 1271591
#> 3     <NA> =IMAGE("https://www.google.com/images/srpr/
#> 4     $A$1                   =ADDRESS(1,1)      NA
#> 5     <NA> =SPARKLINE(R[-4]C[-4]:R[0]C[-4])  NA

```

	A	B	C	D	E	F
1	integer	number_formatted	number_rounded	character	formula	formula_formatted
2	123456	654,321	1.23	one	Google	3.18E+05
3	345678	12.34%	2.35		1,271,591.00	52.63%
4	234567	1.23E+09	3.46	three	Google	0.22
5		3 1/7	4.57	four	\$A\$1	123,456.00
6	567890	\$0.36	5.68	five		317,898

```

cf <- gs_read_cellfeed(gs_ff())
cf %>%
  filter(row > 1, col == 6) %>%
  select(value, input_value, numeric_value) %>%
  readr::type_convert()
#> #> <tibble [5 x 3]>
#> #>   value           input_value numeric_value
#> #>   <chr>          <chr>        <dbl>
#> 1  3.18E+05 =average(R[0]C[-5]:R[4]C[-5])  3.178978e+05
#> 2  52.63%      =R[-1]C[-5]/R[1]C[-5]  5.263144e-01
#> 3  0.22        =R[-2]C[-5]/R[2]C[-5]  2.173942e-01
#> 4  123,456.00 =min(R[-3]C[-5]:R[1]C[-5])  1.234560e+05
#> 5  317,898    =average(R2C1:R6C1)    3.178978e+05

```