

# Reproducible quantitative proteomics analyses using Jupyter notebooks and R

Phillip Wilmarth

Proteomics Core, OHSU

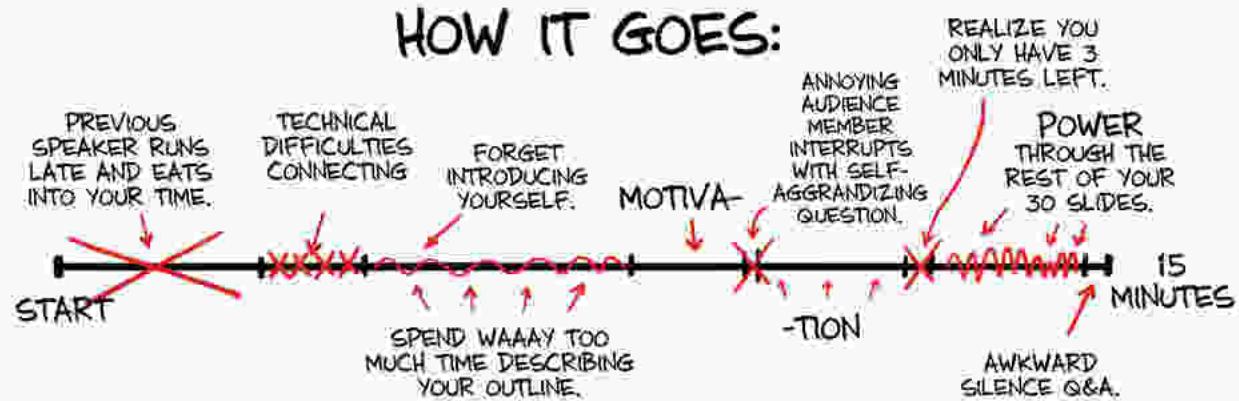
2018 Cascadia Proteomics Symposium

# YOUR CONFERENCE PRESENTATION

## HOW YOU PLANNED IT:



## HOW IT GOES:



Somers J. “The scientific paper is obsolete.” The Atlantic. 2018 Apr 5.



# Atlantic article says notebooks can replace papers – so what are notebooks?

- Notebooks are web server applications
- You work on notebooks with web browsers
- Notebooks format text, run code, and show output
- Jupyter notebooks or R Markdown (RStudio) + knitr

# Installing Jupyter notebooks

- Part of scientific python distributions
  - [www.anaconda.com](http://www.anaconda.com)
- Extensive documentation online
  - [www.jupyter.org/documentation](http://www.jupyter.org/documentation)
- Supports Python by default
- R support needs separate installation
  - <https://irkernel.github.io/installation/>

```
Last login: Sun Jul 15 11:49:24 on ttys001  
[P-J-Imac-3:~ pwilmart$ cd /Users/pwilmart/Box\ Sync/Github_misc/TMT_analysis_examples/Dilution_series  
[P-J-Imac-3:Dilution_series pwilmart$ jupyter notebook  
[I 08:34:43.692 NotebookApp] JupyterLab beta preview extension loaded from /Users/pwilmart/anaconda/lib/python3.6/site-packages/jupyterlab  
[I 08:34:43.692 NotebookApp] JupyterLab application directory  
[I 08:34:43.700 NotebookApp] Serving notebooks from local direc  
tion_series  
[I 08:34:43.700 NotebookApp] 0 active kernels  
[I 08:34:43.700 NotebookApp] The Jupyter Notebook is running at  
[I 08:34:43.700 NotebookApp] http://localhost:8888/?token=fb1  
[I 08:34:43.700 NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).
```

## Launch Jupyter notebook from command line

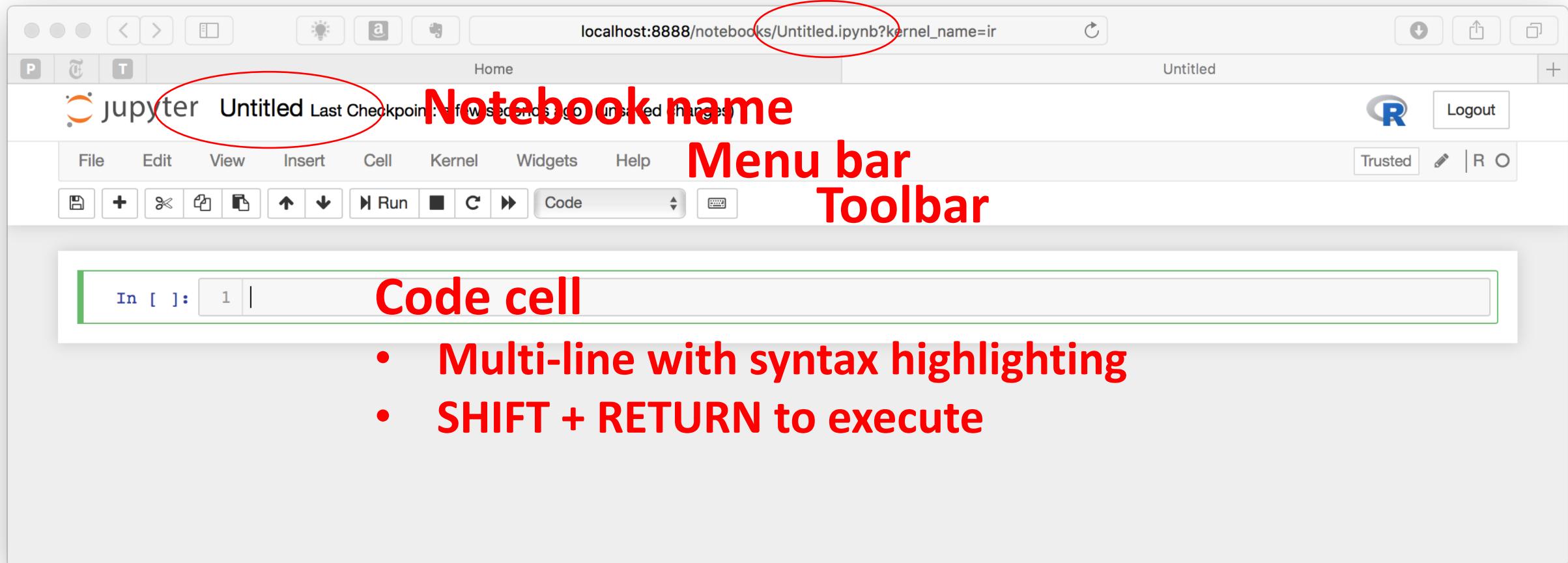
The screenshot shows the Jupyter Notebook interface running in a terminal window. The title bar indicates the current directory is 'Dilution\_series' and the window size is 132x12. The terminal output shows the command-line steps to launch Jupyter Notebook, followed by its startup logs.

The Jupyter interface includes a toolbar with various icons, a header bar with 'localhost:8888/tree', and a main area displaying a file tree. The file tree shows several files and folders:

Name	Last Modified	File size
MAN1353_peptides_proteins.ipynb	a month ago	3.43 MB
grouped_peptide_summary_TMT_8.csv	a month ago	2.55 MB
grouped_protein_summary_TMT_8.csv	a month ago	306 kB
MAN1353_peptides_proteins.html	a month ago	3.65 MB
MAN1353_peptides_proteins.r	a month ago	8.34 kB
psm_tmt.csv	a month ago	11.5 MB
README.md	a month ago	3.52 kB

Red circles highlight specific elements: the terminal command 'jupyter notebook', the 'localhost:8888/tree' URL in the header, the 'New' button in the toolbar, and the 'MAN1353\_peptides\_proteins.ipynb' file in the file tree.

# Notebook anatomy



The screenshot shows a Jupyter Notebook interface with the following details:

- Header:** The title bar displays "localhost:8888/notebooks/MAN1353\_peptides\_proteins.ipynb".
- Toolbar:** The top menu includes File, Edit, View, Insert, Cell, Kernel, Widgets, Help, and a dropdown set to "Markdown".
- Cell Content:** A code cell contains the following text:

```
1 # Dilution Series and Data Aggregation
2 ## Phil Wilmarth
3 ### OHSU PSR Core
4 ### February, 2018
5 <br>
6
7 # 1. Compare PSM, peptide, and protein data scatter
8 This is a complex mouse brain mixture, digested, split into 6 aliquots, labeled with 6 TMT reagents, and mixed in
9 a series of dilutions. The relative volume amounts were in these proportions: 25 : 20 : 15 : 10 : 5 : 2.5. The
10 biological sample was from the Gail Mandel lab at OHSU.
11 The data were generated on a Thermo Fusion using the SPS MS3 method (McAlister 2014). The peptide identifications
12 and reporter ion intensities were generated using an in-house pipeline with Comet (Eng 2013) as the search engine.
13 The PAW pipeline (Wilmarth 2009) does a direct processing of the RAW files to extract MS2 scans for Comet and the
14 reporter ions from the corresponding MS3 scans (Chambers 2012). Comet results are filtered using a target/decoy
15 strategy with accurate mass conditional score histograms. Basic parsimony logic is used to generate protein and
16 peptide reports.
17 An additional protein grouping algorithm is also employed to group protein families having large fractions of
18 identical common peptides. The final list of proteins is used as the context for shared and unique peptide
determinations. Only unique peptides in that context are used in the TMT quantification. All unique PSM reporter
ions are summed into protein intensity totals for each channel. We will compare some properties of the reporter
ions at the PSM level, the peptide level (combined copies and charge states), and the protein level.
```
- Annotations:** A red oval highlights the first four lines of the code cell, which are bolded. To the right of the cell, the text "header levels" is written in red. A red circle highlights the word "text" in the 12th line of the code cell, and the text "indented text" is written in red below it.

localhost:8888/notebooks/MAN1353\_peptides\_proteins.ipynb

jupyter MAN1353\_peptides\_proteins Last Checkpoint: 06/15/2018 (unsaved changes)

Home MAN1353\_peptides\_proteins Logout

File Edit View Insert Cell Kernel Widgets Help Trusted | R O

Dilution Series and Data Aggregation

Phil Wilmarth

OHSU PSR Core

February, 2018

## 1. Compare PSM, peptide, and protein data scatter

This is a complex mouse brain mixture, digested, split into 6 aliquots, labeled with 6 TMT reagents, and mixed in a series of dilutions. The relative volume amounts were in these proportions: 25 : 20 : 15 : 10 : 5 : 2.5. The biological sample was from the Gail Mandel lab at OHSU.

The data were generated on a Thermo Fusion using the SPS MS3 method (McAlister 2014). The peptide identifications and reporter ion intensities were generated using an in-house pipeline with Comet (Eng 2013) as the search engine. The PAW pipeline (Wilmarth 2009) does a direct processing of the RAW files to extract MS2 scans for Comet and the reporter ions from the corresponding MS3 scans (Chambers 2012). Comet results are filtered using a target/decoy strategy with accurate mass conditional score histograms. Basic parsimony logic is used to generate protein and peptide reports.

An additional protein grouping algorithm is also employed to group protein families having large fractions of identical common peptides. The final list of proteins is used as the context for shared and unique peptide determinations. Only unique peptides in that context are used in the TMT quantification. All unique PSM reporter ions are summed into protein intensity totals for each channel. We will compare some properties of the reporter ions at the PSM level, the peptide level (combined copies and charge states), and the protein level.

Chambers, M.C., Maclean, B., Burke, R., Amodei, D., Ruderman, D.L., Neumann, S., Gatto, L., Fischer, B., Pratt, B., Egertson, J. and Hoff, K., 2012. A cross-platform toolkit for mass spectrometry and proteomics. *Nature biotechnology*, 30(10), p.918.

Eng, J.K., Jahan, T.A. and Hoopmann, M.R., 2013. Comet: an open-source MS/MS sequence database search tool. *Proteomics*, 13(1), pp.22-24.

# Notebooks are a linear collection of cells

GitHub, Inc. GitHub.com/pwilmart/TMT\_analysis\_examples/blob/master/Dil.R

```
3rd Qu.: 573798 3rd Qu.: 456530 3rd Qu.: 349698 3rd Qu.: 241573
Max. :30487165 Max. :23870707 Max. :18641957 Max. :12991318
e_5 f_2.5
Min. : 50 Min. : 50
1st Qu.: 9015 1st Qu.: 4778
Median : 32801 Median : 18037
Mean : 149827 Mean : 82819
3rd Qu.: 114691 3rd Qu.: 63668
Max. :6002579 Max. :3401632
```

### Prepare the data for plotting

We will use the average of the 6 intensities as the x-axis and plot each dilution series (the different reporter ion channels) against that. We will use different colors for each dilution series channel and overlay 6 scatter plots in one figure. We will show the scatter plots with linear axis scales, so we can see the details. We will add a linear fit line for each series.

```
In [67]: # create an average vector for the x-axis
psms$ref <- rowMeans(psms)
peptides$ref <- rowMeans(peptides)
proteins$ref <- rowMeans(proteins)

# we can simplify plotting if we put data in long form (tidy data)
gpmss <- gather(psms, key = dilution, value = intensity, a_25:f_2.5)
gpeptides <- gather(peptides, key = dilution, value = intensity, a_25:f_2.5)
gproteins <- gather(proteins, key = dilution, value = intensity, a_25:f_2.5)

# check some things
head(gproteins)
```

ref	dilution	intensity
94285.463	a_25	188305.8
50322.493	a_25	106780.4
6875.984	a_25	13065.7
86251.341	a_25	160570.5
532073.475	a_25	1044356.3
1297868.915	a_25	2533375.2

### Same thing with some transformations for MA plots

We will also make data frames for MA style plots. Those plots have transformed axes (log2 ratios on y-axis, and log10 average intensity on x-axis). We will do the average intensity in log10 scale (that is easier to mentally "unlog" to get back intensities). We will do the ratios in log2 scale since that is easier for estimating fold changes.

```
In [68]: # make frames for MA style plots
log_psms <- log2(psms[1:6] / psms$ref)
log_psms$ref <- log10(psms$ref)
log_peptides <- log2(peptides[1:6] / peptides$ref)
log_peptides$ref <- log10(peptides$ref)
log_proteins <- log2(proteins[1:6] / proteins$ref)
log_proteins$ref <- log10(proteins$ref)

# also tidy the log data frames
glog_psmss <- gather(log_psms, key = dilution, value = log_ratios, a_25:f_2.5)
glog_peptides <- gather(log_peptides, key = dilution, value = log_ratios, a_25:f_2.5)
glog_proteins <- gather(log_proteins, key = dilution, value = log_ratios, a_25:f_2.5)

# compute the ratios of each dilution channel to the reference
# save as log values for horizontal lines in the MA plots
calc_ratios <- colMeans(psms)
calc_ratios <- log2(calc_ratios[1:6] / calc_ratios[7])

# check some things
round(calc_ratios, 2)
```

a_25	0.92
b_20	0.61
c_15	0.22
d_10	-0.28
e_5	-1.38
f_2.5	-2.22

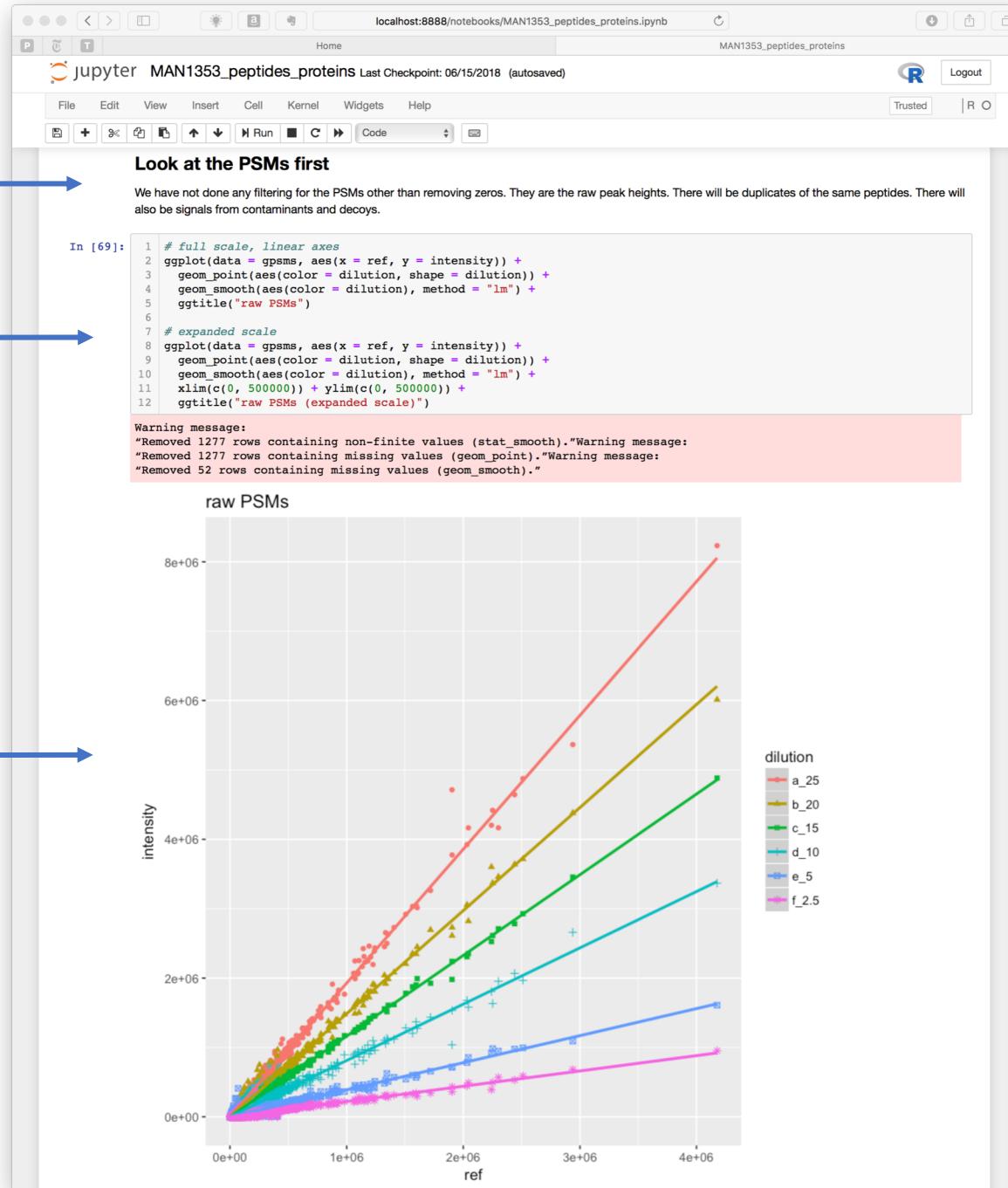
# Pictures are worth a thousand words

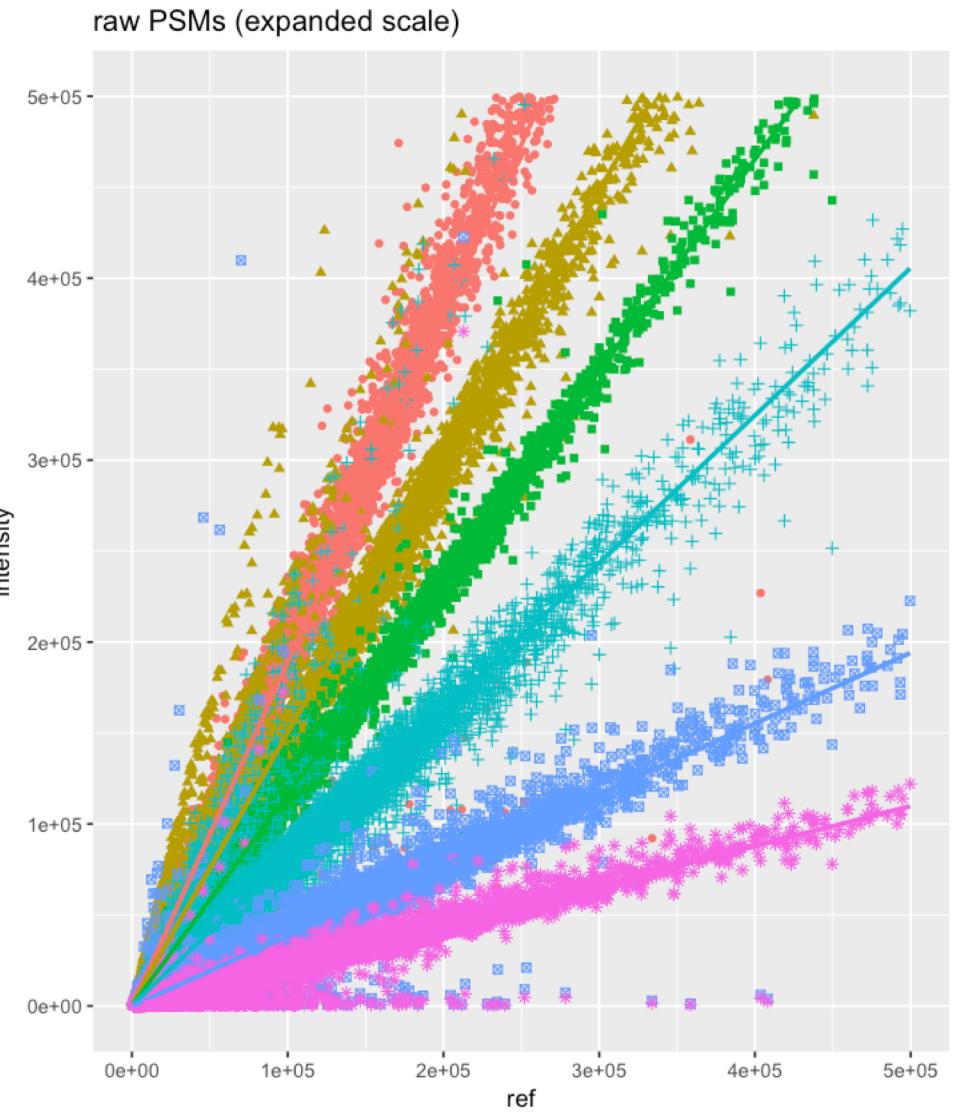
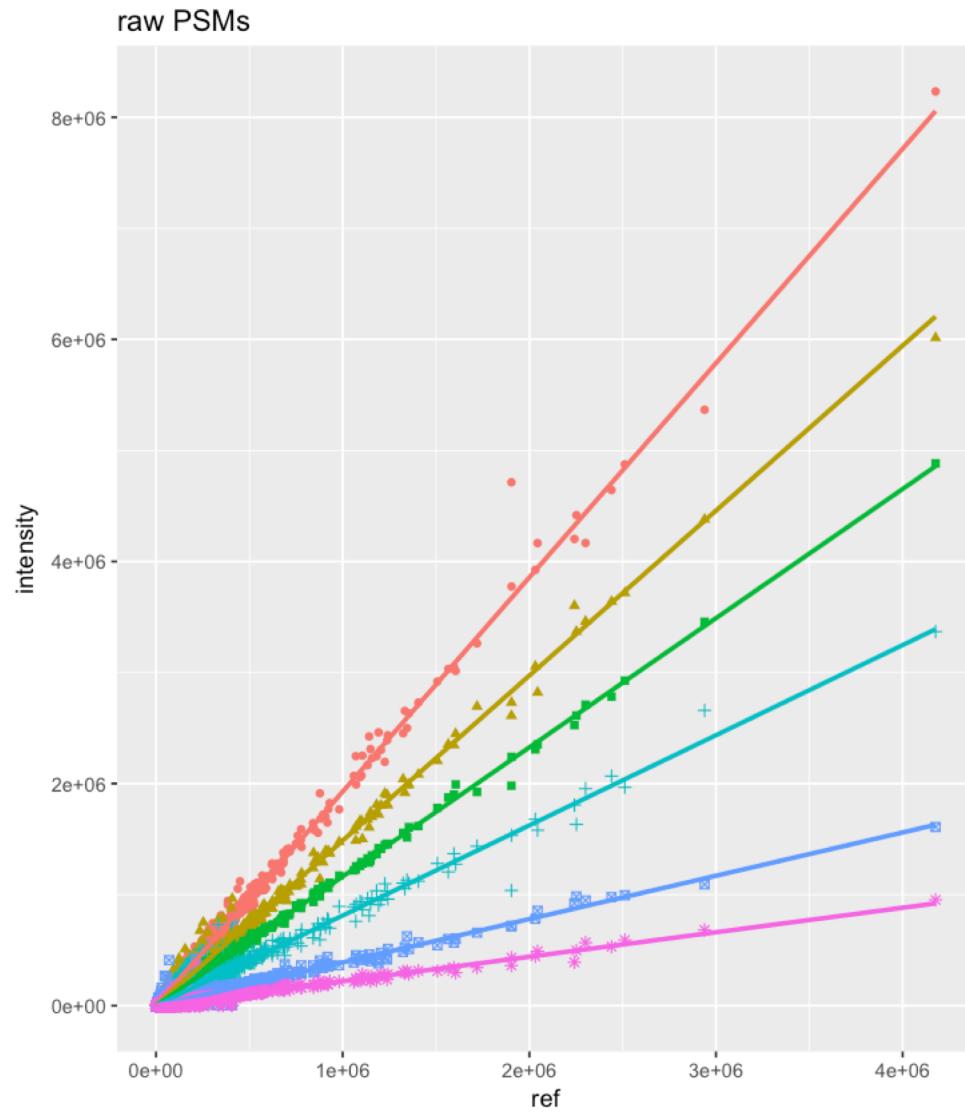
- R and ggplot2 are very powerful for data visualization
- TMT labeling of same digest in 6 dilutions
  - 25 : 20 : 15 : 10 : 5 : 2.5
- Comparisons at PSM, peptide, protein levels
- [https://github.com/pwilmart/TMT\\_analysis\\_examples](https://github.com/pwilmart/TMT_analysis_examples)

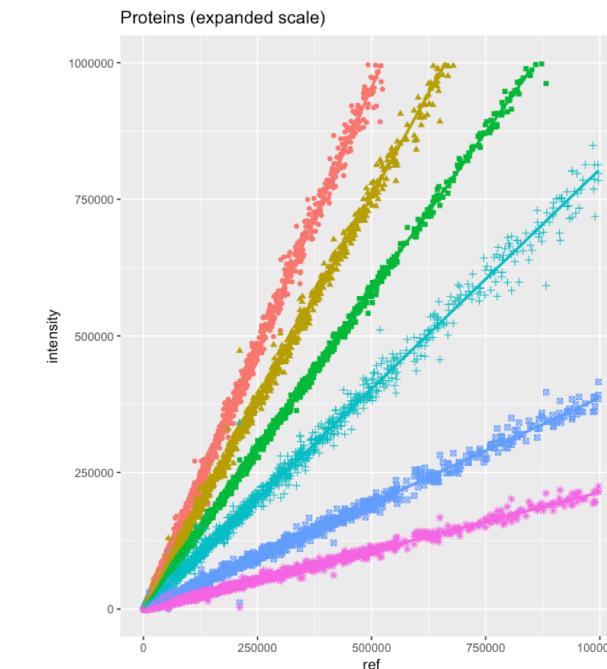
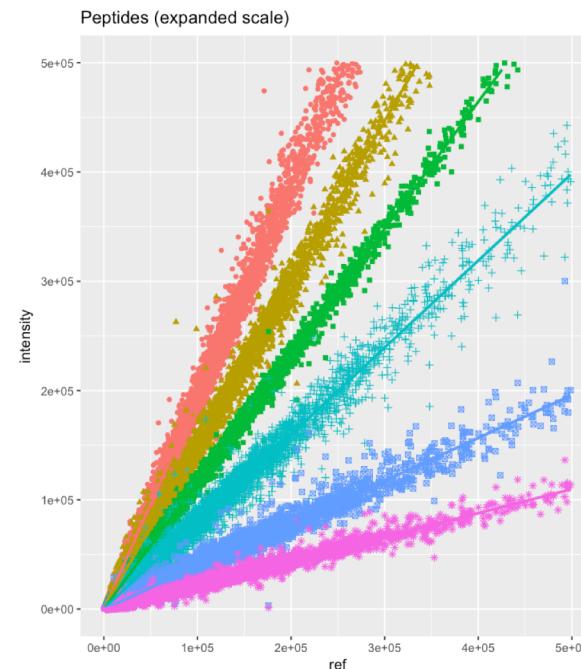
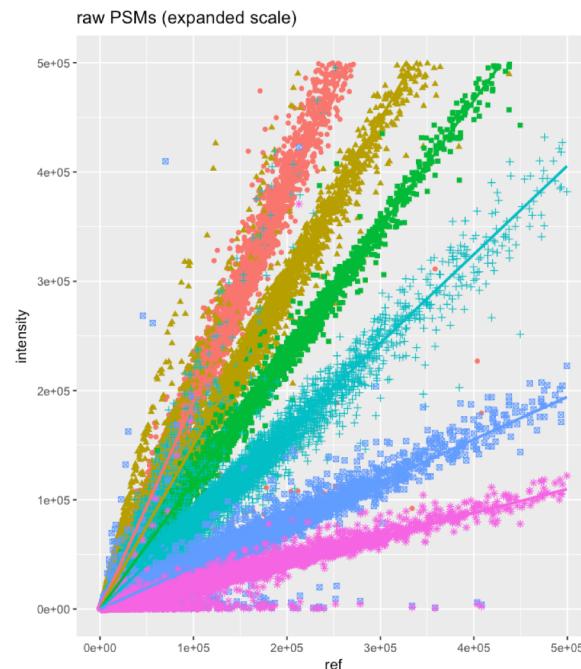
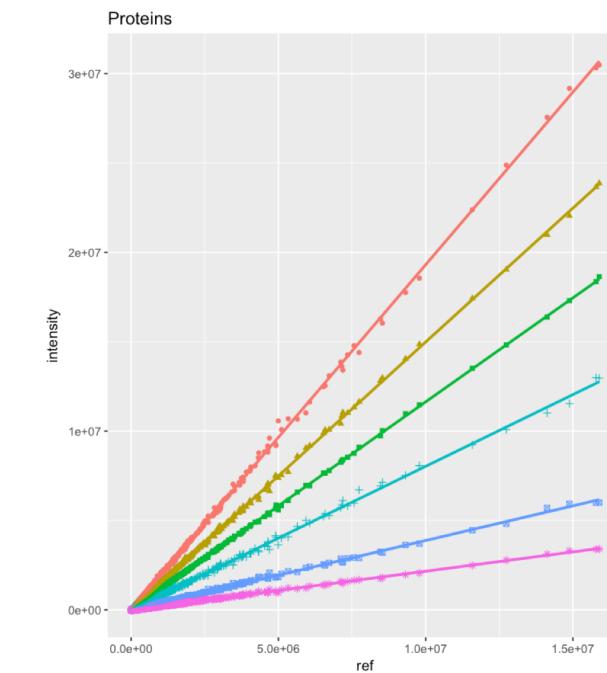
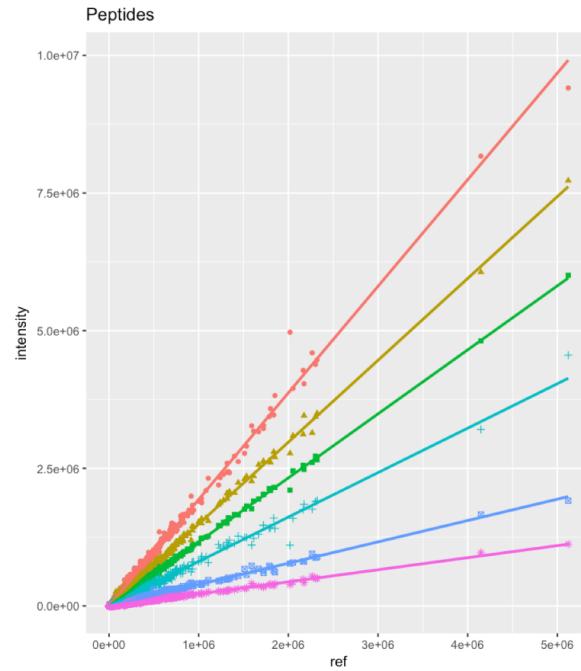
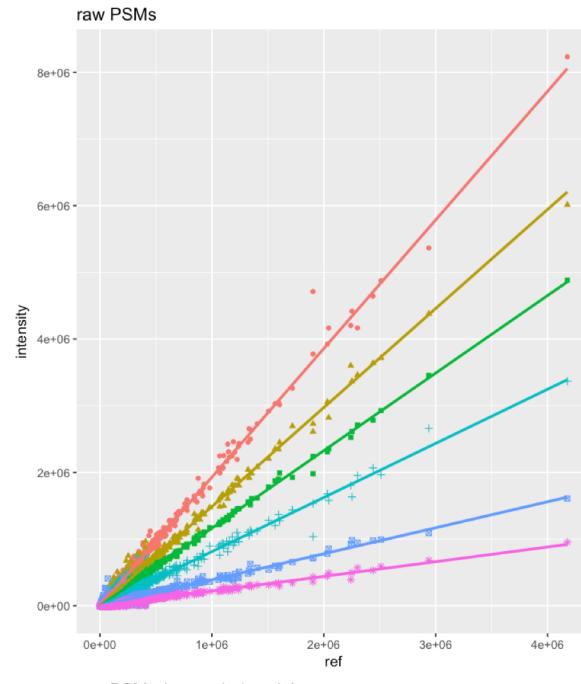
# Markdown

# Code

# Output

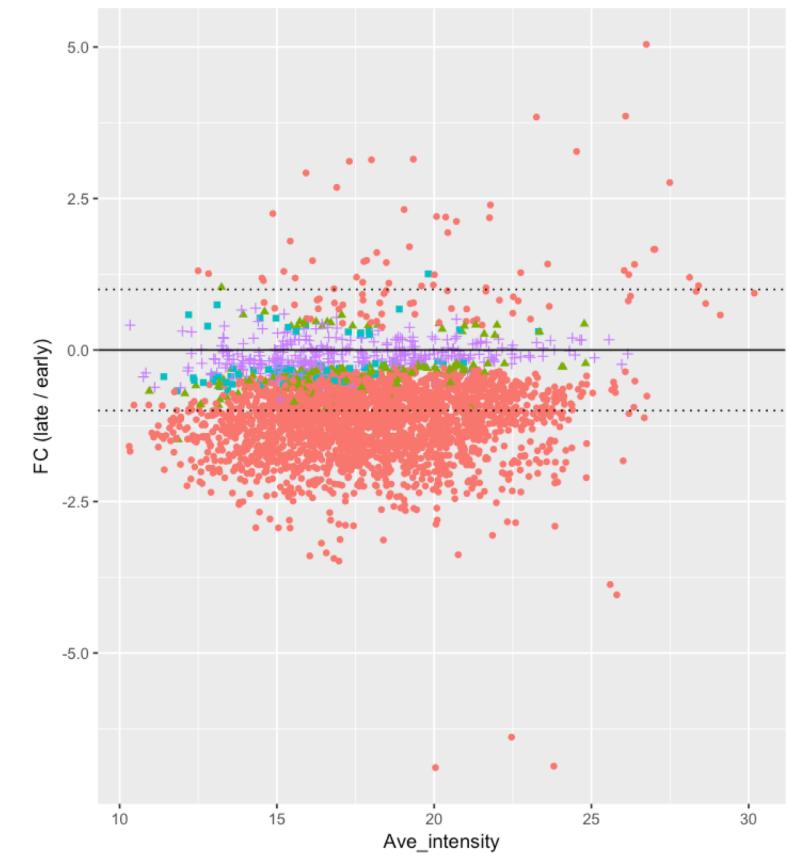




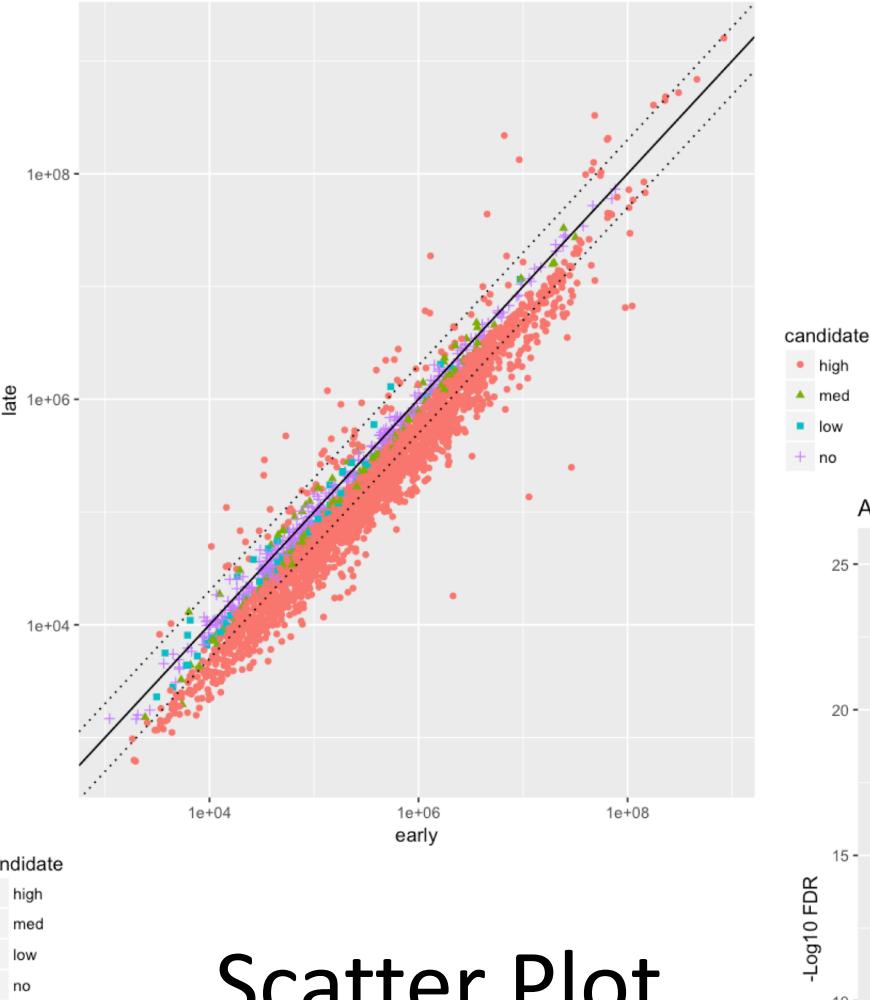


# MA Plot

After IRS early vs late (MA plot)



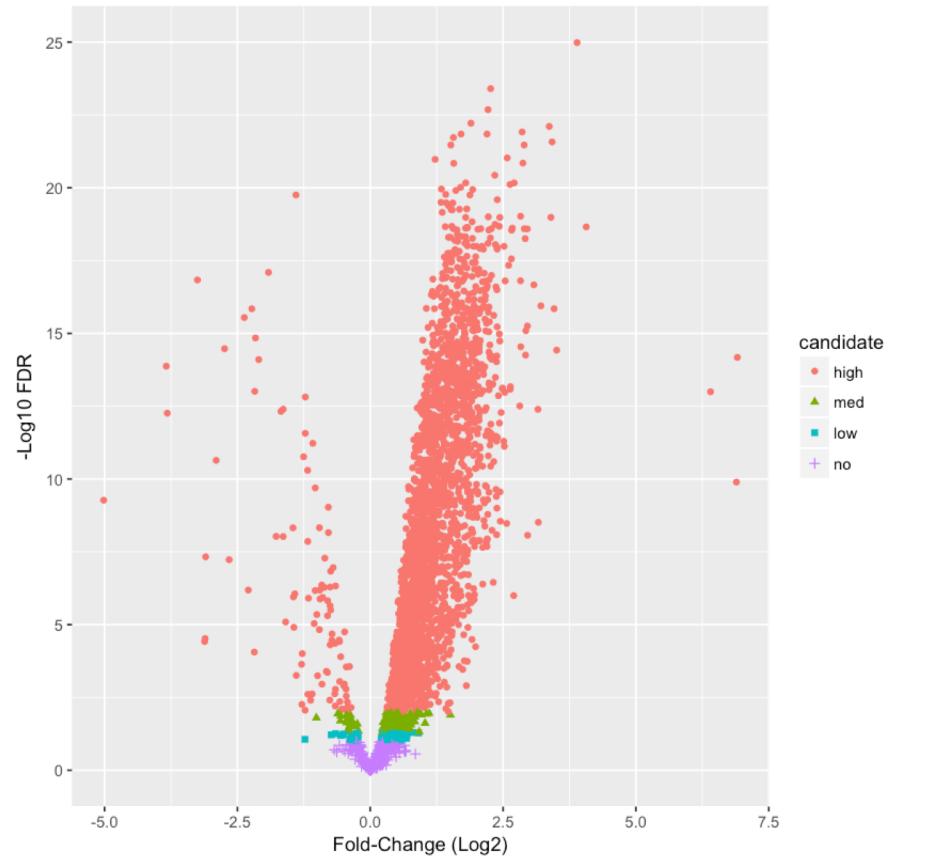
After IRS early vs late (scatter plot)



# Scatter Plot

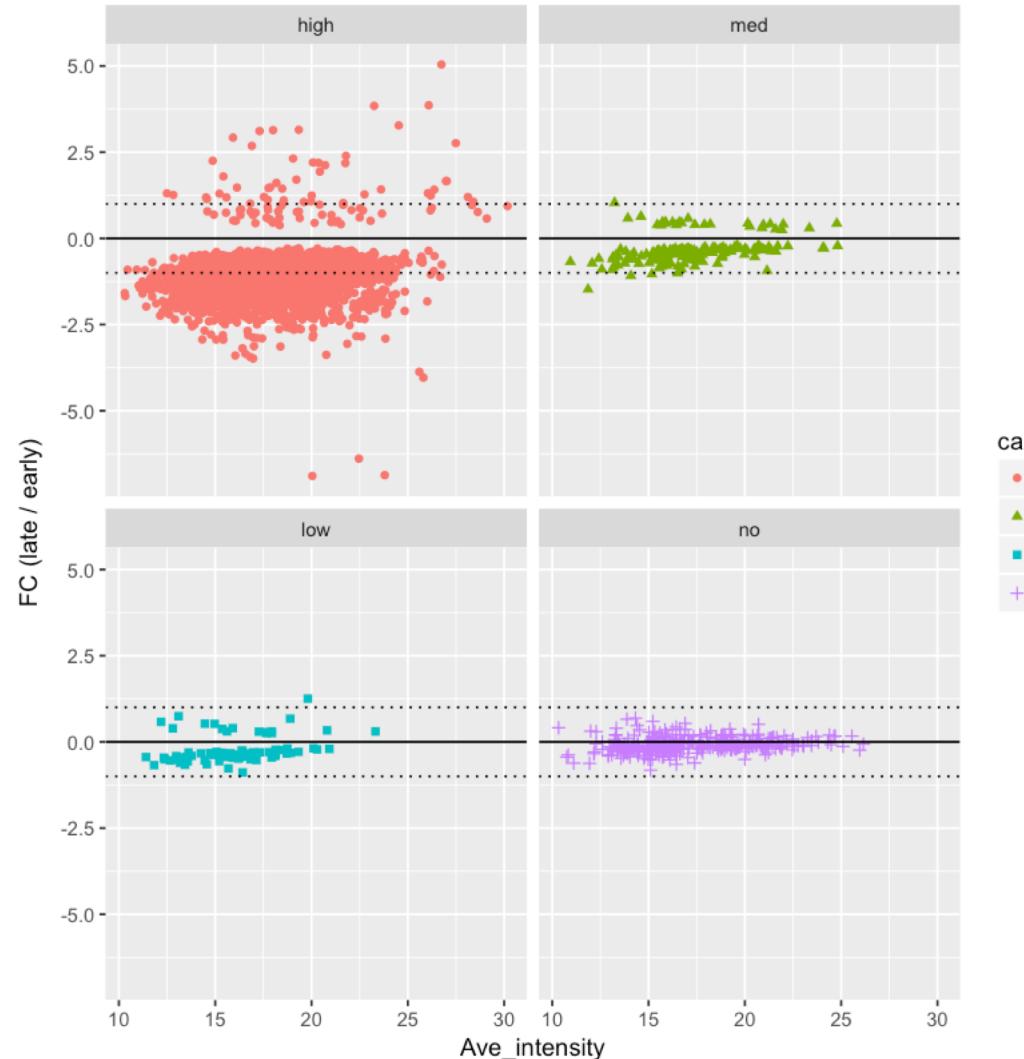
# Volcano Plot

After IRS Volcano Plot

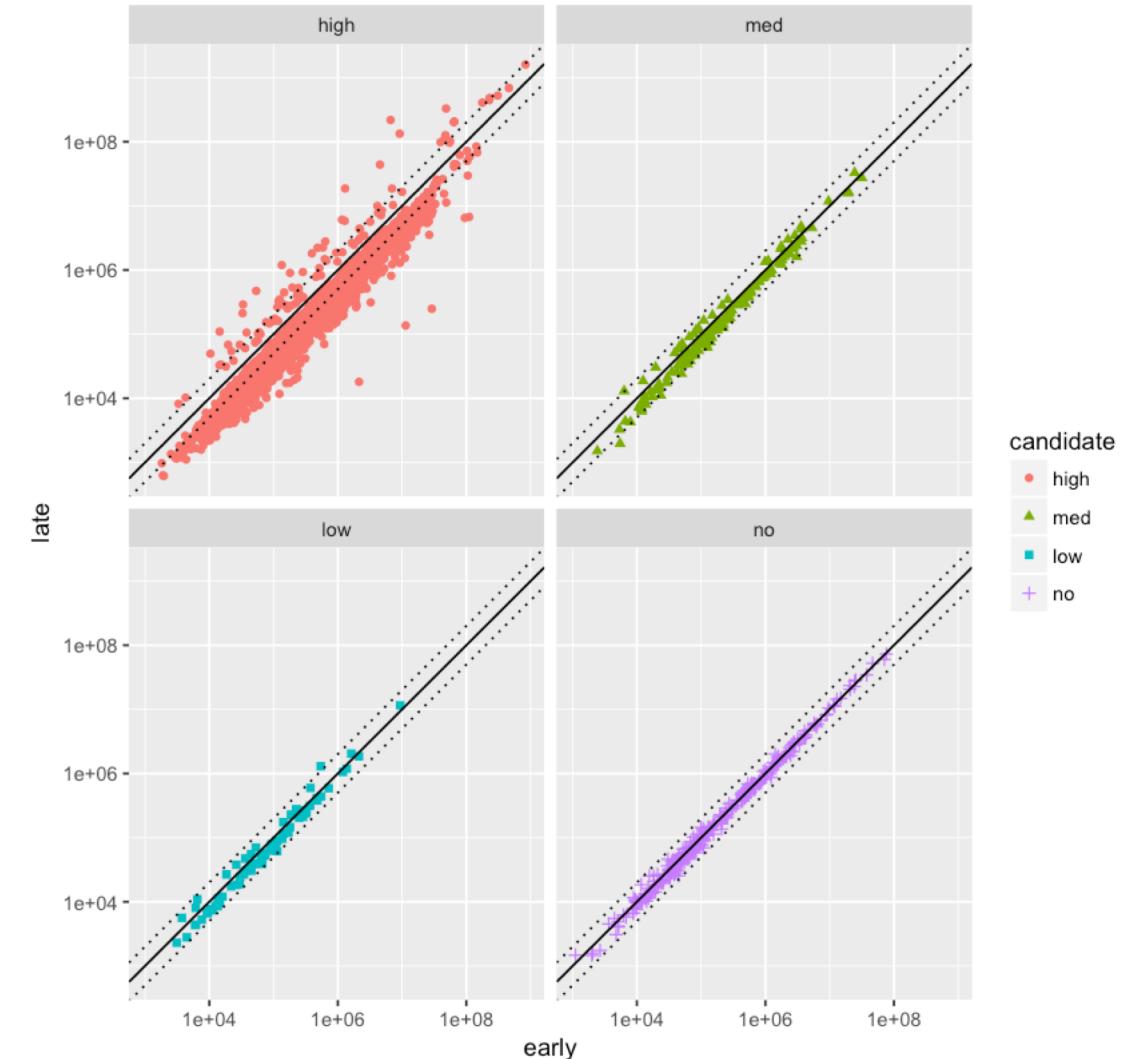


# Facet plots – separate candidates by significance

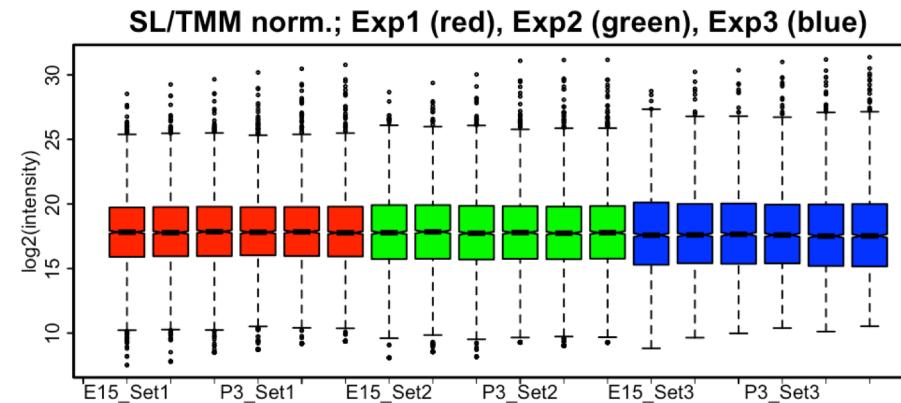
After IRS MA plots, separated by candidate



IRS scatter plots, separated by candidate

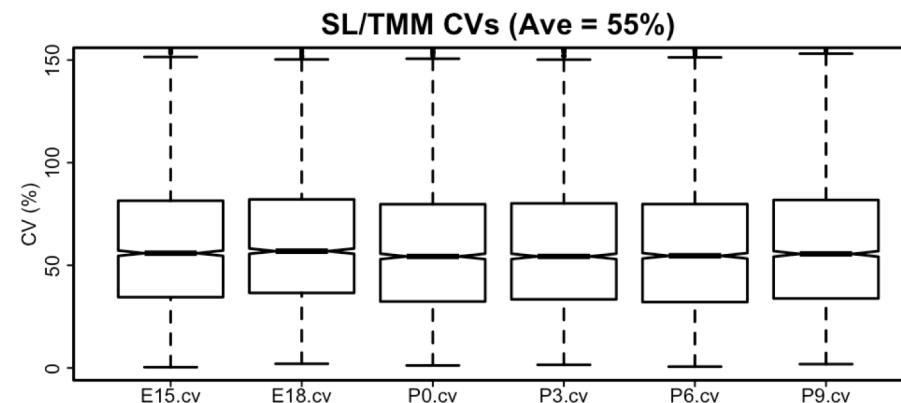
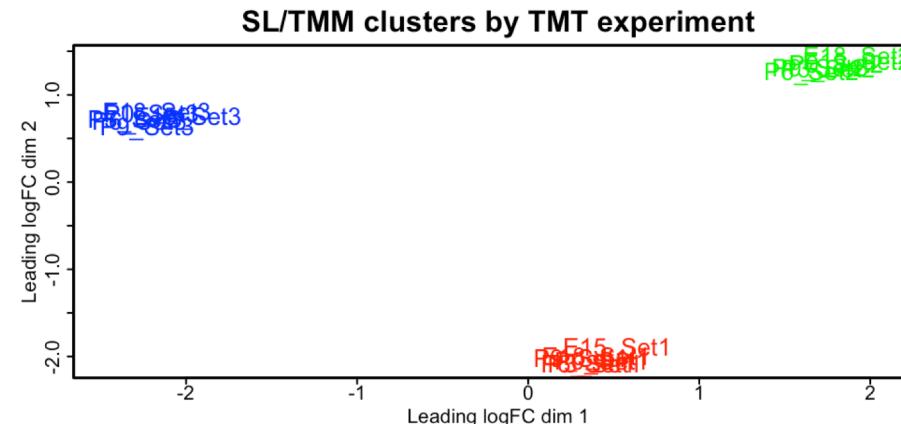


# Without IRS



Mouse lens: 6 developmental times  
biological replicates in separate TMT experiments

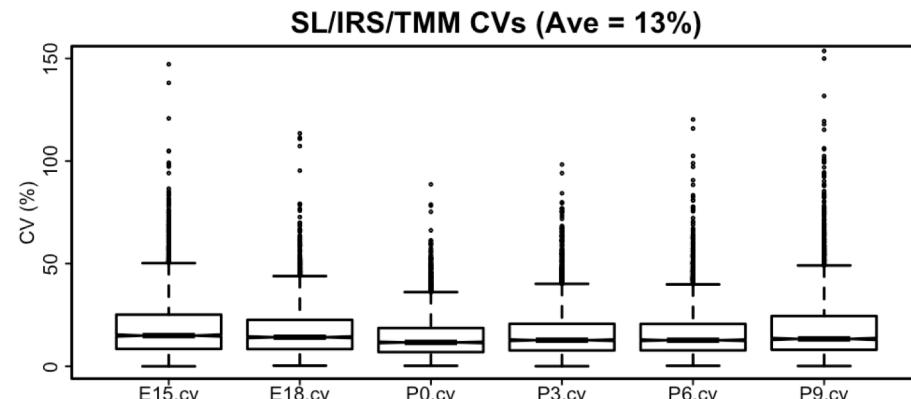
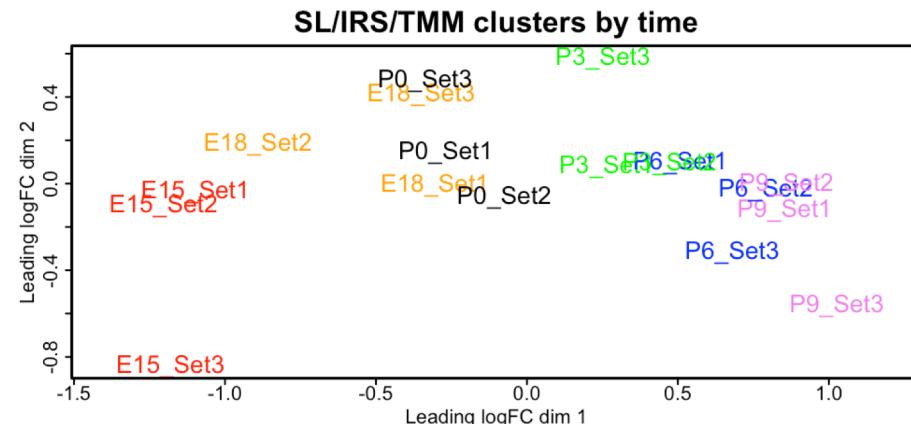
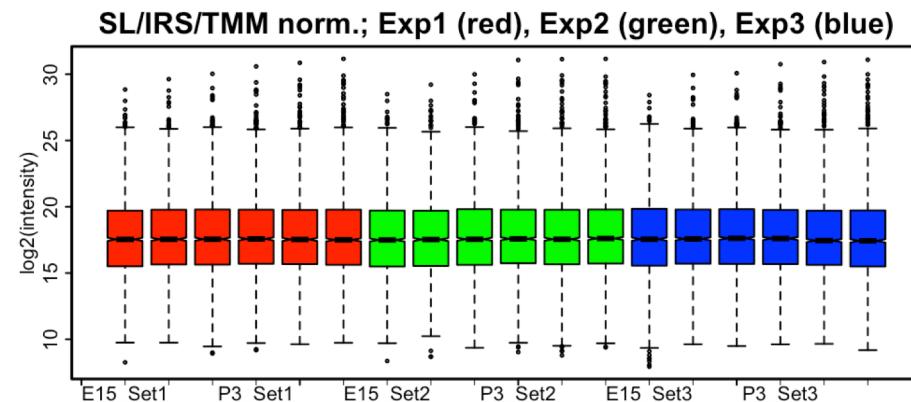
Single factor normalizations seem to work



They don't actually work

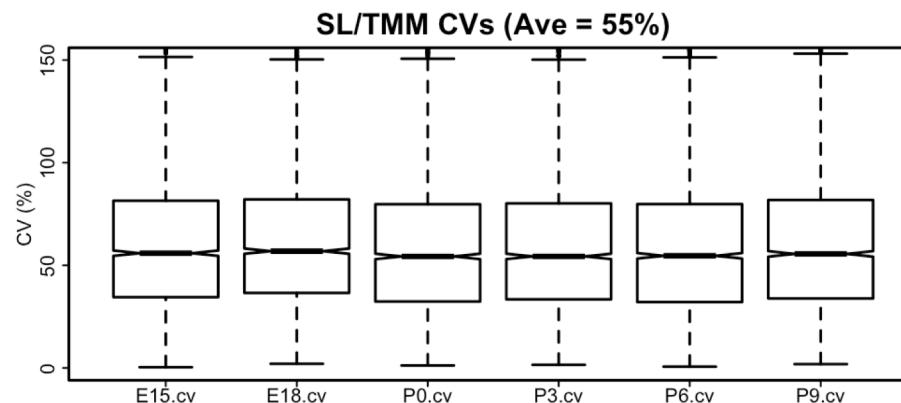
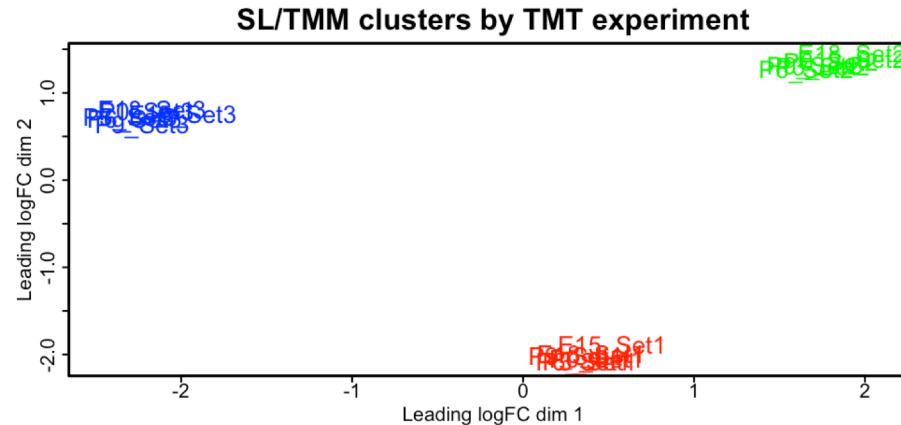
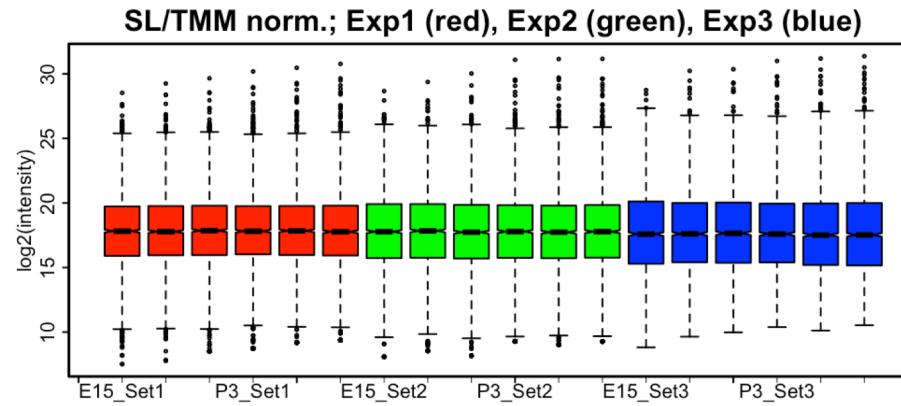
Internal Reference Scaling:  
corrects random MS2  
sampling that affects  
reporter ion intensities  
between scans.

**With IRS**

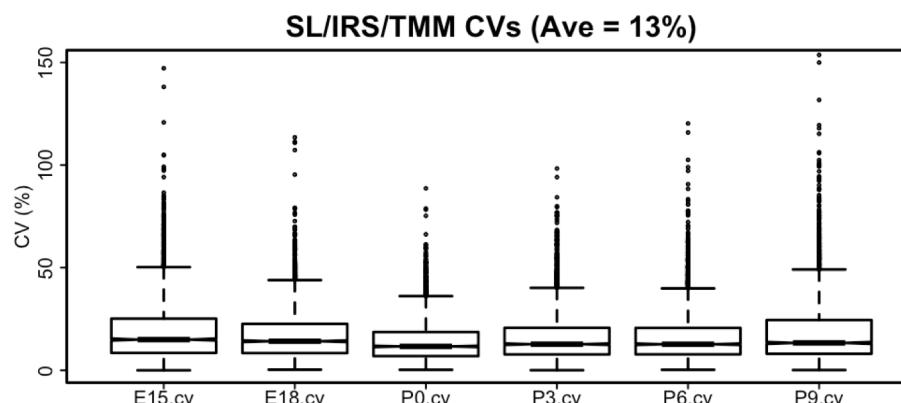
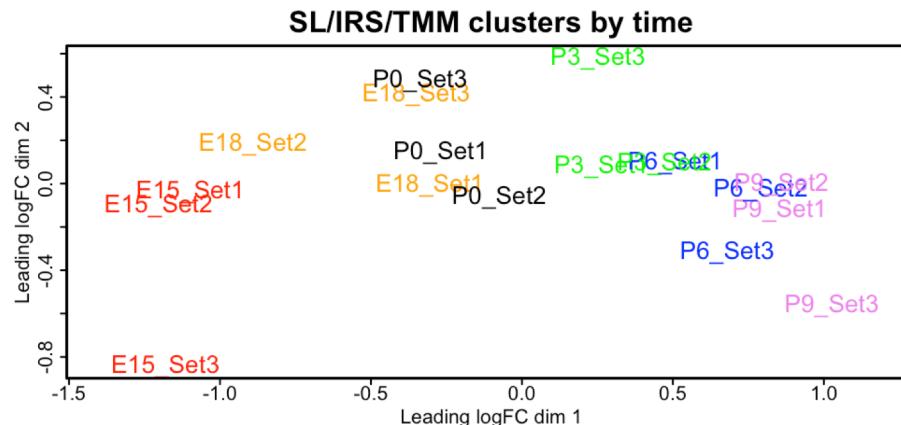
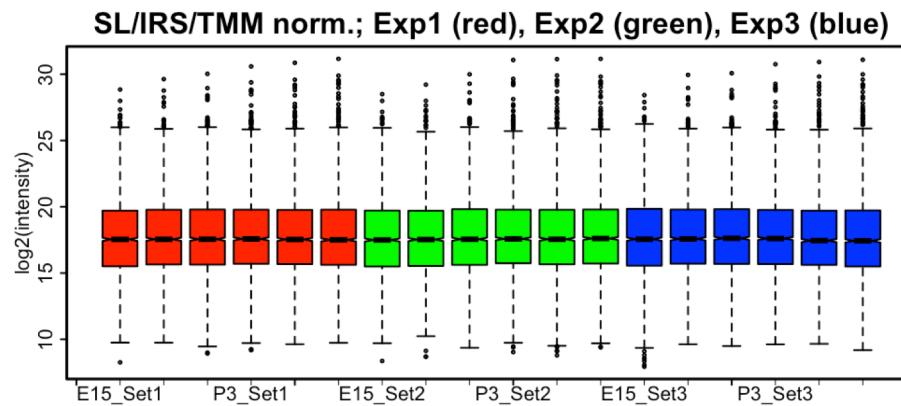


Plubell, et al., MCP, 2017, v16, p873-890.

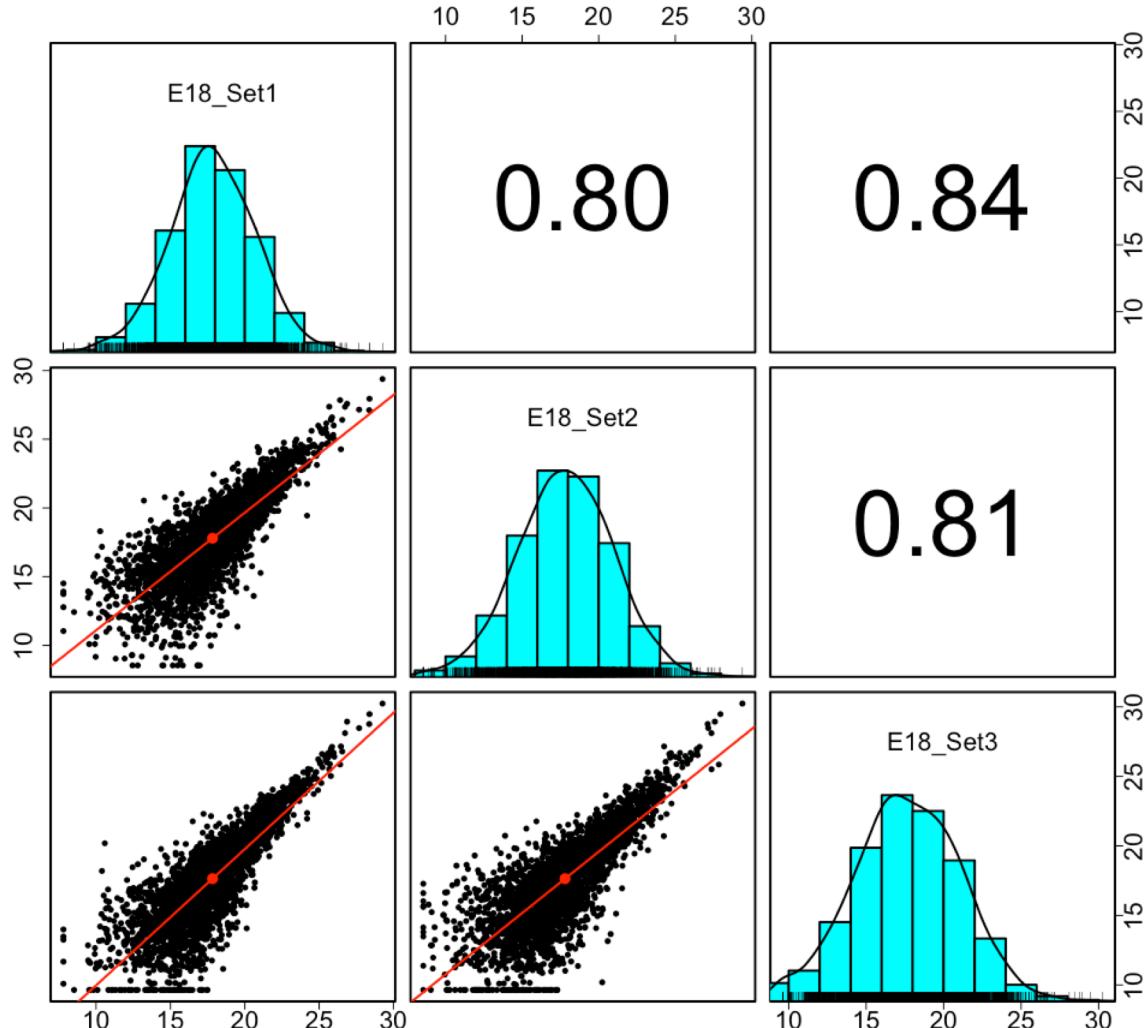
# Without IRS



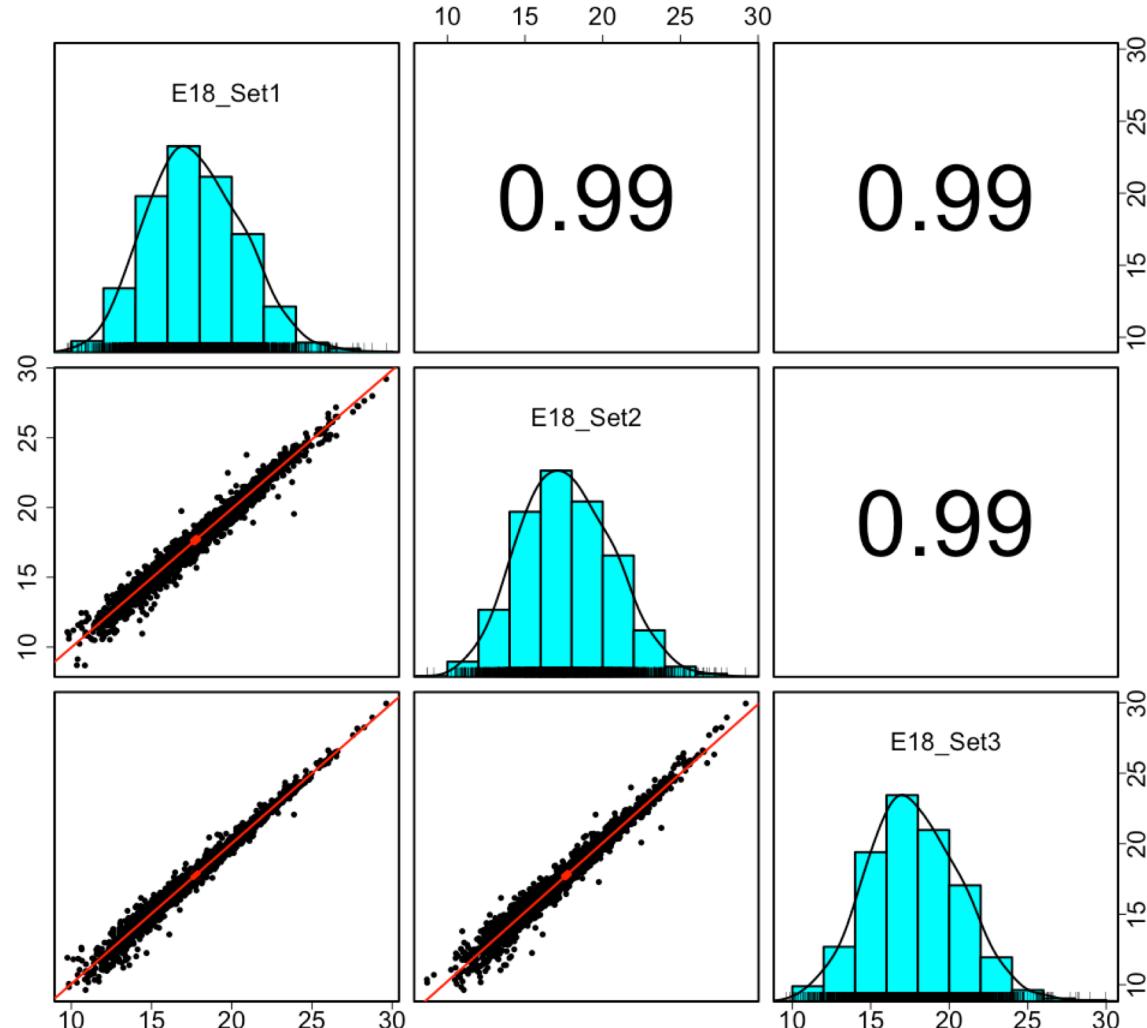
# With IRS



# Without IRS



# With IRS



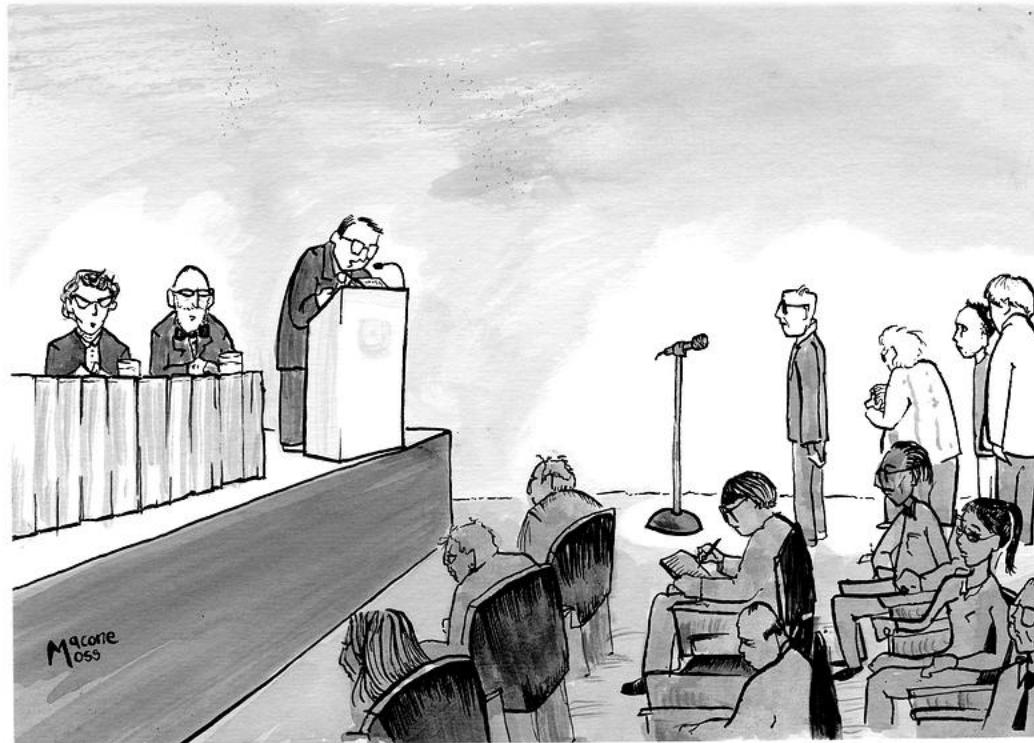
# Atlantic article says notebooks can replace traditional scientific papers...

- Notebooks don't work for all content
- Notebooks make great Supplemental Materials
- Notebooks are logical choices for:
  - QC and sanity checks
  - Normalizations, transformations, correlations
  - Statistical testing and modeling

# Disclaimers

- Scientific python and Jupyter notebooks are big topics
- R is confusing and there is a lot to learn
  - Traditional (old) R
  - RStudio and modern R
  - Statistics
  - Bioconductor
- Git/Github also extensive
- “data science” field is interesting: Hadley Wickham, Wes McKinney, Jake VanderPlas (UW)

See: <https://github.com/pwilmart>



*"We'd now like to open the floor to shorter speeches disguised as questions."*

<https://goo.gl/images/x2PZXP> (Steve Macone, New Yorker, October 18th, 2010)