

Problem: Excel changes my input data

Month: JAN to DEC

a number contains the letter E or number contains more characters than can be displayed

a number contains leading zeros

Note: Copy from HTML page: use paste special -> Unicode text

You will Learn:

- General tips and tricks in navigation, filtering, and selection
 - Keep input data types
 - Keep headers visible as you scroll
 - Keyboard navigation shortcuts
 - Quickly identify patterns within a column
 - Split items within single column
 - Highlight cells by their values
- · Using formulas to work with data
- Merge multiple tables with vlookup
- Summarize data with PivotTable

Whitehead Institute

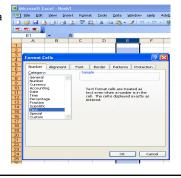
2

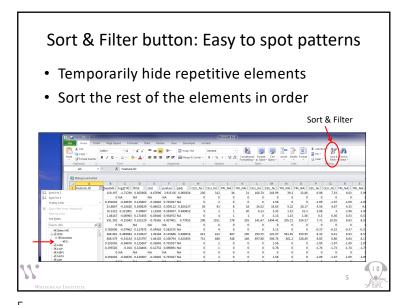
Solution: Change formatting to "Text"

1. Open file with plain text

Before you populate the data into Excel, select the cells that will contain the gene symbols, right-click and choose Format

Cells >Text > OKCopy and paste data

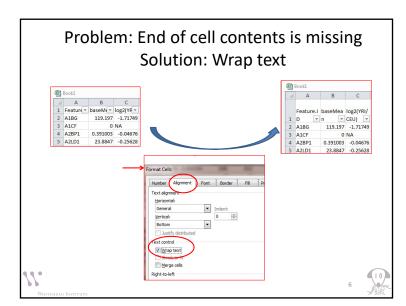


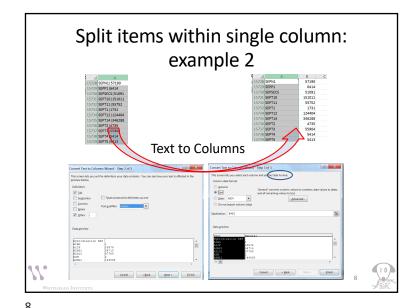


Split items within single column:
example 1

Text to Columns

Text to Colu





Merge multiple columns =CONCATENATE(A1, "-", B1, "-", C1) 7 TCG4-48-407B 4007 7 TCGA-A8-A08B A007 TCGA A030 A00Z TCGA A04T 01A 21R A034 TCGA A034 A04V 01A 21R TCGA-A2-A04\ 10 TCGA A041 01A 11R A034 11 TCGA A046 01A 21R A034 12 TCGA A04A 01A 21R A034 13 TCGA A04C 01A 21R 4034 A04D 01A 21R A034

Keyboard navigation shortcuts PC MAC First cell in the active column Control/Ctrl Up arrow Command Last cell in the active column Control/Ctrl Down arrow Command First cell in the active row Control/Ctrl Left arrow Last cell in the active row Control/Ctrl Right arrow Command The very first cell (A1) Control/Ctrl Home Control/Ctrl Control/Ctrl **Bottom right corner** Control/Ctrl End plus

11

Choose View > Freeze Panes -> Freeze top Row To unfreeze frozen rows, choose View > Freeze Panes > Unfreeze Panes Wicrosoft Excel John Zoom 1 New Arran Vindow All View Side by Side John Polition Workspare Selection Window All Vindow All View Side by Side John Workspare Selection Window All Vindow All View Side by Side John Workspare Selection Window All Vindow All View Side by Side John Workspare Selection Window All Vindow All View Side by Side John Workspare Selection Window All Vindow All View Side by Side John Workspare Selection Window All Vindow All View Side by Side John Workspare Selection Window All View Arran Vindow All View Side by Side John Workspare Selection Window All View Arran Vindow All View Side by Side John Workspare Selection Window All View Arran Vindow All View Side by Side John Workspare Selection Workspare Selection Window All View Side by Side John Workspare Selection Workspare Selec

10

Keyboard selection shortcuts

- Selecting a large range of cells can be a pain. Here's one trick that could help.
 - To start the selection of a large matrix, start by selecting the bottom right cell of the desired matrix (Ctrl + End)
 - Use "Ctrl + Shift + up arrow" to select everything above the original cell.
 - Use "Ctrl + Shift + left arrow" to select everything to the left of the original cell.
 - Use "Shift + down arrow" to move down one row (to avoid the header line, if present).
 - Use "Shift + right arrow" to move to the right one column (to avoid any ID column(s), if present)

Whitehead Institute

The Cancer Genome Atlas (TCGA)

- NCI / NHGRI project to molecularly characterize
 - Over 20k primary cancer samples
 - Matched normal samples
 - 33 cancer types
- Genomic, epigenomic, transcriptomic, and proteomic assays
- · Data we'll use for exercises:
 - Small subset of BRCA RNA-seq data from TCGA
 - Unit: normalized gene expression levels

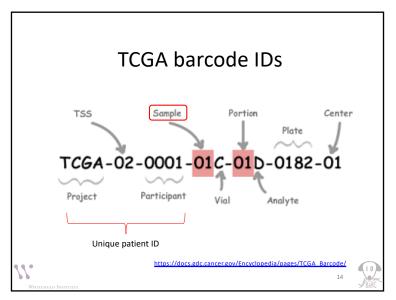
13



Whitehead Inst

13

TCGA Sample codes Code Short Letter Code 01 Primary Solid Tumor 02 Recurrent Solid Tumor Primary Blood Derived Cancer -03 Peripheral Blood Tumor (01-09) Recurrent Blood Derived Cancer - TRBM Bone Marrow Additional - New Primary TAP Metastatio TM 07 Additional Metastatic TAM 08 Human Tumor Original Cells THOC Primary Blood Derived Cancer 09 TBM Bone Marrow 10 Blood Derived Normal Solid Tissue Normal NT Buccal Cell Normal NRC 13 EBV Immortalized Normal NFR\ 14 Bone Marrow Normal NBM 15 sample type 15 15SH sample type 16 16SH Control Analyte CELLC Complete sample codes can be found at:

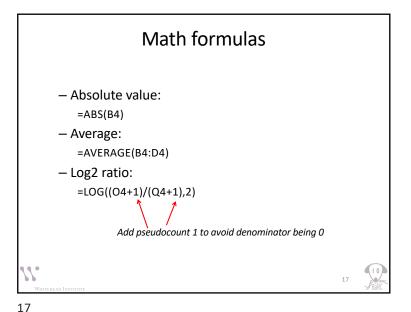


14

Ex.1 (TCGA_data.xlsx)

- Freeze the portion of the sheet BRCA_from_TCGA, so you can always see the column names (TCGA barcodes) and gene names no matter how you scroll.
 - hint: View ->Freeze panes / Freeze Top Row
- Split the TCGA barcodes and check the sample codes. How many sample types?
 - hint: copy and paste barcodes to another sheet with Paste special ->
 Transpose, then split barcodes with Data -> Text to Columns -> View sample codes by turning on Sort & Filter button under Home panel)
- Replace sample barcodes with meaningful names:
 - 01 -> Primary
 - 06 -> Metastasis
 - 11 -> Normal
- Add sample codes as additional row to top of the BRCA_from_TCGA sheet.
 - hint: Copy -> Paste special -> Transpose
- · Sort the samples by sample codes
 - hint: Data -> Sort -> Options -> Sort left to right

RaRC

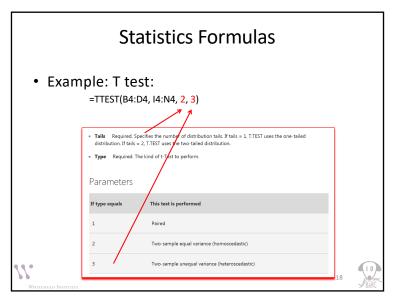


Logical formulas

• Examples:

- If, else
=IF(U4<0.05, 1, 0)
False

- and or
=IF (AND (T4<0.05, R4>=1), "yes", "no")



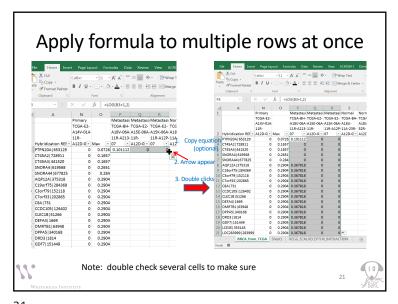
18

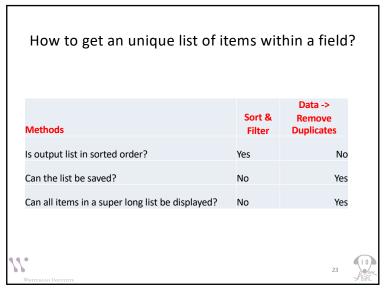
20

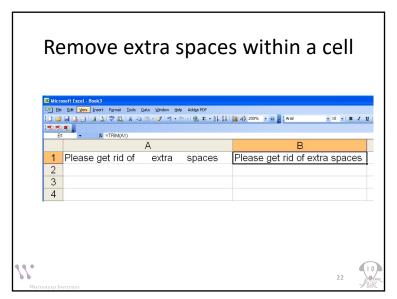


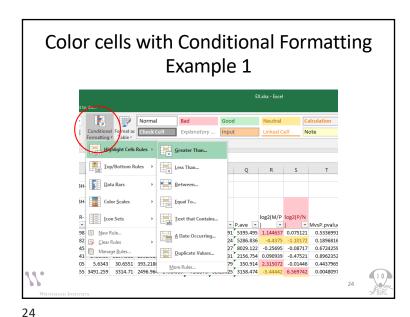
- Begin with equal sign: =
- Dollar sign (\$): absolute cell references
- Convert formula to real value:
 - Copy → Paste Special → Values

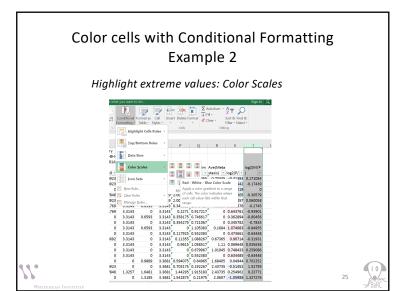
20











- Drop genes with no expression in all samples
 - Hint: Create another column and name it as "Max" -> use formula to calculate maximum value for each gene -> Sort data by "Max"

Ex.2 (TCGA data.xlsx)

- 2) How do you easily find out where the non-zero value start? Hint:
 - a) Convert values in "Max" column from formula format to value with Copy -> Paste Special -> Values
 - After sorting, with Sort & Filter button to find out the value right above 0 -> Search for this value under "Max"
- Compare primary tumor vs normal, and metastatic cancer vs primary tumor:
 - . Convert values to log2 format:
 - To avoid dividing by 0, you can add 1 to each value
 - Start to get the log2 expression for the top gene in the first patient sample, use our tricks (see slides) to fill out the rest of the cells
 - Calculate log2 average for each gene in normal samples, and do the same for primary tumor and metastasis samples.
 - Calculate log2 fold changes (log2FC) between primary tumor and normal, and between metastatic cancer and primary tumor, and highlight the log2FC with Conditional Formatting
 - Calculate p-values with ttest, delete the genes with no valid pvalues, find out FDR adjusted pvalue with the BaRC web tool (http://bioinfo.wi.mit.edu/bio/tools/fdr/), and highlight FDR adjusted p-values < 0.05 with Conditional Formatting.
- We consider a gene is "differentially expressed" if meeting the cutoffs (FDR adjusted p-value < 0.05 + at least 2 fold difference). Add two columns named as Sig(PvsN) and "Sig(MvsP)", fill in with "y" and "n" to indicate if a gene is "differentially expressed" or not.

Hint: with "if" plus "and" formula

Barc

Color cells with Conditional Formatting Example 3 Highlight extreme values: Color Scales -> More Rules View ACROBAT Developer ⊞ - | ♂ - A - | ■ ■ | ■ ■ | Merge & Center - \$ - % • 100 → 00 Conditional Form ■ Format only cells that contain ► Format only top or bottom ranked values Format only values that are above or below average ➤ Format only unique or duplicate values og2(M/P 0.384316 -0.34068 MAC: 0.093457 0.077152 0.64 Format values that rank in the -0.30015 -1.06402 0.54 Excel version 16: ✓ % of the selected range

✓ Top/bottom 0.205713 0.109701 It's "Top/Bottom Rules" -0.97894 -1.33834

-0.53919 2.732735 1.403292 -2.04853

0.448461 -0.88095

0.008764 -0.00214

26

Lookup and retrieve data Microsoft Excel - Microarray_Analysis_data 2_SOLUTION [Read-Only] 🕙 Elle Edit <u>V</u>iew Insert Format Iools Da<mark>ta W</mark>indow Help Ado<u>b</u>e PDF 🗋 📸 🖟 💪 (△) ኞ 🛝 | ¾ 🖎 (△) | ♥ - ♥ - | № Σ - ½↓ ¼↓ (∭ - № - 0) 💂 Arial =VLOOKUP(\$E2,\$A\$2:\$C\$489,2,FALSE) 2 100_g_at 5 603044225 7 114704126 11 39205046 1160 at 9.368257073 4 1052 s at 6.935502676 7.634905024 5 1063 s at 7 177596401 9.29804293 6.356422623 3.554644254 3.59319203 59528921 8.04178479 1 1194_g_at 5 287191623 1.79435423 2 1232_s_at 3.68424095 5.938234238 14 1264_at 15 1271_g_at 4 638605837 3 465669364 7.456925352 7.1958078 Too complicated? Try BaRC Submatrix Selector http://bioinfo.wi.mit.edu/bio/tools/submatrix selector/

27

Look up values in data list: vlookup

- VLOOKUP(lookup_value,table_array,col_index_num,ran ge_lookup)
 - lookup_value: The value to search in the first column of the table array
 - table_array: The table to search (containing the value to search for in the first column)
 - · Note: use "\$" to fix the table
 - col_index_num: the column number from which the matching value is returned
 - range_lookup:
 - True/omitted, an exact/approximate match
 - False: an exact match

Dollar sign (\$): absolute cell references

29



WHITEHEAD IN

29

In this example, it grouped the values by gene and calculated the average

Ex.3 (TCGA_data.xlsx)

 Extract gene symbols from the cells under column "Hybridization REF". How do you check gene symbols to make sure none of them is converted to date format?

Copy cells in this column and Paste them next to the rightest column -> use "Text to Columns" to split cells

- Use vlookup to get the expression levels of a list of genes, such as genes in the KEGG_ECM_RECEPTOR_INTERACTION gene set from Broad's MSigDB.
- 3. Are any of these genes "differentially expressed" as defined by
 - 1) adjusted p-value < 0.05
 - e) absolute log2 fold change >= 1

Hint: Use "Sort & Filter" button to quickly narrow down the genes

WHITEHEAD INSTITUTE

30

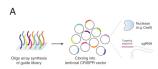


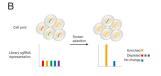
30

32

Pooled CRISPR screening

- Simultaneous targets of a large number of genes or regulatory elements.
- To prevent off-target, library includes multiple single guide RNAs (sgRNAs) for each gene.
- The entire cellular population is subjected to selective pressure.
- The genes driving the phenotype of interest are identified by comparing the relative abundance of corresponding sgRNAs in both fractions.





Anal Biochem. 2017 532:95

32

Ex.4 (CRISPR data.xlsx)

- CRISPR_data.xlsx includes the percentage of guide RNA expression in a subset of a pseudo pooled CRISPR screen experiment. Each gene has 4 single guide RNAs.
- Could you calculate the percentage of expression at gene level with PivotTable? You can summarize expression of the sgRNAs by taking their average percentage.



33

- Next Hot Topics:
 - Enrichment Analysis January 28
- List of Winter/Spring Hot Topics Sessions:
 - http://barc.wi.mit.edu/education/hot_topics/upc oming/

35



References

- Excel Annoyances Curtis Frye
- Excel Hacks O'Reilly
- http://www.exceltip.com (Joseph Rubin)
- http://excel.tips.net (Allen Wyatt)
- Some Excel Basics as well as formula basics:
 - http://barc.wi.mit.edu/education/bioinfo2005/arrays/Exce | help.html
- Zeeberg BR, Riss J, Kane DW, Bussey KJ, Uchio E, Linehan WM, Barrett JC, Weinstein JN. BMC Bioinformatics 5:80 2004

