# Data: X

$n \times p$

individuals    variables ("genes")    "reads"

# RNA-Seq:

Quantification done by "read counts"

$X_{ij}$: $\dfrac{\text{\# of times read } j \text{ was observed}}{\text{in sample } i.}$ Counts

What we have currently:
- Not $X$
- Raw data files for each individual:
  - FASTQ

Step 1:
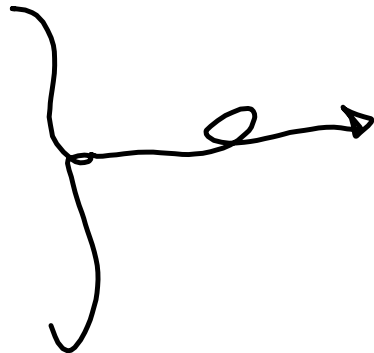- Process the FASTQ files into $X$

⌐ reference R docs I sent.

Step 1b: Create $Z_{n \times 4}$ matrix of "meta data".

# Step 2 :

$X$ : Which of the "genes" differ between

$n \times p$      diseased & healthy, adjusted for

the meta data variables.

Gene 1 : p-value

:               List of "differentially

expressed" genes

Statistical analysis:

Gene-by-gene:

$$y_{ij} = \beta_0 + \beta_1 TX_j + \beta_2 Act_j + \cdots$$

$$H_0 : \beta_1 = 0 \quad \text{vs.} \quad H_a : \beta_1 \neq 0$$

In R: edgeR, DESeq2, limma

Instead of p-values, we'll actually use alternative significance measure called False Discovery Rate (FDR)

Next Steps:
- Create Github user account.
  - I'll create repository

- Do some background reading
  - R/Bioconductor tutorial
  - Khan Academy for basic (2 hrs)
    genetics & cell/molecular biology
  - You guys web conference to brainstorm,
    assign tasks & due dates.

R Resources:

  ▷ R Cookbook

  ▷ Codeschool.com:

    ▷ "Try R" tutorial