
Associating gene expression to neural network derived image features

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 I will be updating this document with the progress of my NIPS 2017 submission.

2 Plan

3 Q1: How much variability is shared between image features and expression?

- 4 1. Calculate PCA of image features and expression
- 5 2. Calculate cross-correlations between the PCs until 99.9% variance is explained.
- 6 3. Use the cross-correlation measurements to derive total fraction of variance explained of
7 images by expression: calculate $\sum_i \lambda_i \sum_j r_{ij}^2$ where λ_i is the i th eigenvalue of the image
8 feature PCA, r_{ij}^2 is the correlation of image feature PC i with expression feature PC j .
- 9 4. Make all of the above one-line-to-run for any combination of patch size, feature layer, tissue
- 10 5. Report the value for each patch size, feature layer, tissue; interpret findings

11 Q2: How much of the shared variability is due to technical confounders?

- 12 1. Regress out known confounders C from the image feature F and expression data E . I.e.,
13 fit models $F \sim C + \epsilon$ and $E \sim C + \epsilon$. Use log scale for Ischemic time, monitor all scatter
14 plots for confounders that are included in the model that the linear model assumptions hold.
15 Decide whether to retain full model of all confounders x all genes/features, retain only
16 significant associations, or do forward feature selection.
- 17 2. Repeat steps Q1:1-5 above using the residuals of this model.
- 18 3. Create a scatter plot of total shared variance vs. technical shared variance [again, all tissues
19 etc.]
- 20 4. Interpret findings in context of many tissues, scales, feature layers

21 Q3: What genes share expression variability to visual features, and at which scales?

- 22 1. Using residuals of technical confounders, fit linear model of expression vs. image feature -
23 calculate effect sizes and p-values
- 24 2. Apply multiple testing correction to derive a list of image feature associated genes.
- 25 3. Sort genes based on variance explained by each image feature, perform GSEA or gProfiler
26 ranked gene list analysis to test whether particular aspects of biology are enriched. Interpret
- 27 4. Assess number of associations [total and per-feature] for different tissues, feature scales.
28 Assess overlap of findings from gene list analyses. Interpret.

29 **Q3A: What genes are influenced by technical confounders?**

30 1. For each confounder, calculate the fraction of variance it explains for each gene expression
31 level.

32 2. Sort genes based on variance explained by confounder, perform GSEA or gProfiler ranked
33 gene list analysis to test whether particular aspects of biology are enriched.]

34 **Once these done**

- 35 • Q4: What do the image features represent?
- 36 • Q5: Is there a genetic basis to the image features?
- 37 • Q6: Are there gene expression levels that cause image feature changes, or vice versa?
- 38 • Q7: Are there image features that are associated with clinical annotations? Are any causal?

39 **Methods**

40 **Assessing shared variability between image features and expression**

41 We calculate principle components that capture 99.9% of variation in the image features and expres-
42 sion features respectively. We repeat this after fitting a linear model containing known technical
43 confounders to both the image features and expression, and regressing out these effects.

44 We choose the top 2000 transcripts with mean expression greater than one, and look for associations
45 with the image features across tissue samples. For image features generated from Lung tissue, we
46 plot a QQ-plot with

47 **Results**

48 Variability in expression explains a substantial amount of variability in the image features (Figure
49 ???). Furthermore, almost 50% of the variability in expression, and almost 40% of the variability in
50 the image features, is explained by 51 known technical factors (Figure ???)

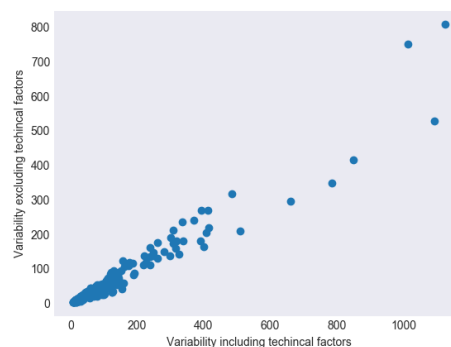


Figure 1: Sample figure caption.

Table 2: Enriched Ontology terms for image feature 3

p-value	Ontology	Genes
0.0498	BRONCHIECTASIS WITH OR WITHOUT ELEVATED SWEAT CHLORIDE 3; BESC3;;CYSTIC FIBROSIS-LIKE SYNDROME	SCNN1G