
Explaining gene expression variation to neural network derived image features

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 How does gene expression influence visual characteristic of tissues? Can the
2 variation in expression be explaining using data extracting from these images
3 alone? To what extent is this possible? In this work, I am to answer these questions,
4 by employing neural networks on a dataset of biomedical images annotated with
5 genotype and gene expression data.

6 1 Introduction

7 Biomedical images are routinely used by doctors in pathology to diagnose diseases. For example,
8 lung biopsies are taken from patients when grading different levels of lung cancers. Specifically for
9 lung cancer, a great deal of work has focussed on the identification of specific tissue characteristics
10 which are indicative of metastasis. Indeed, this has been the topic of many supervised machine
11 learning competition and approaches. We can now accurately identify metastatic regions from high
12 resolution histopathology lung images, with the state of the art being achieve by Convolutional Neural
13 Networks (CNNs).

14 What other types of information can be estimated from these high resolution images? We explore this
15 direction using a datasets of.....

16 2 Results

17 We find that across 10 different tissue types

18 Plan

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

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

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

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

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

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

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

- 46 1. For each confounder, calculate the fraction of variance it explains for each gene expression
47 level.
- 48 2. Sort genes based on variance explained by confounder, perform GSEA or gProfiler ranked
49 gene list analysis to test whether particular aspects of biology are enriched.]

50 **Once these done**

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

55 **Methods**

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

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

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

63 **Results**

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

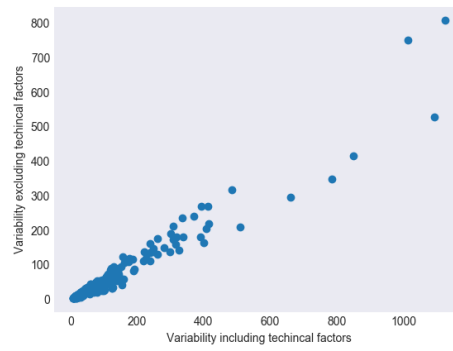
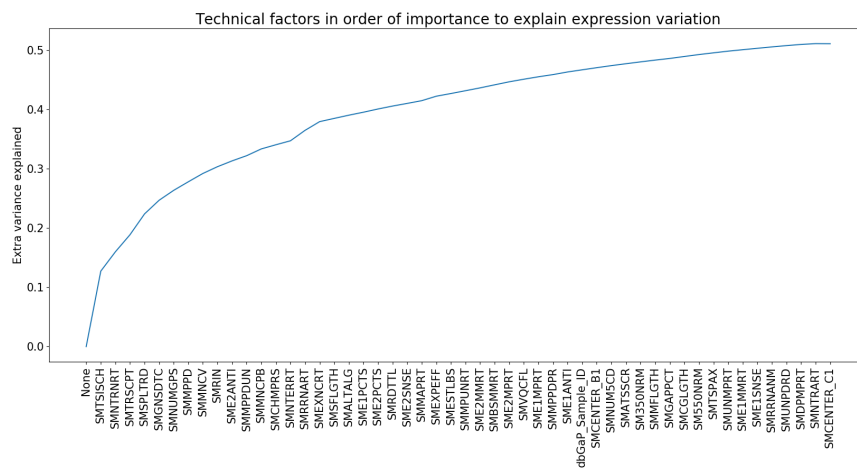


Figure 1: Sample figure caption.



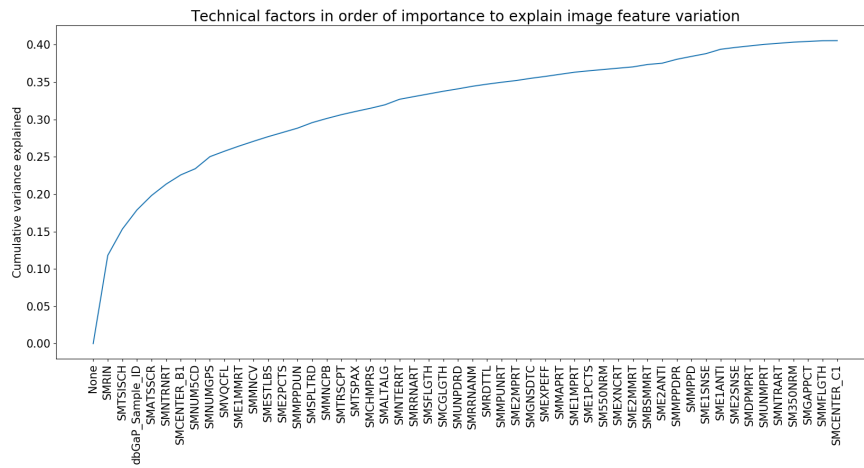


Figure 3: Sample figure caption.

Table 1: Enriched Ontology terms for image feature 501

p-value	Ontology	Genes
1.21e-15	lamellar body	LAMP3,CTSH,SFTPA1,NAPSA,SFTPD
6.87e-09	Defective ABCA3 causes pulmonary surfactant metabolism dysfunction type 3 (SMDP3)	SFTPA1,SFTPD,SFTPB,SFTPA2

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