

# DataRemix: a universal data transformation for optimal inference from gene expression datasets

Weiguang Mao<sup>1,2</sup>, Ryan Hausler<sup>2</sup> and Maria Chikina<sup>1,2\*</sup>

## Abstract

RNAseq technology provides an unprecedented power in the assessment of the transcription abundance and can be used to perform a variety of downstream tasks such as inference of gene-correlation network and eQTL discovery. However, raw gene expression values have to be normalized for nuisance biological variation and technical covariates, and different normalization strategies can lead to dramatically different results in the downstream study. Here we present a simple three-parameter transformation, DataRemix, which can greatly improve the biological utility of gene expression datasets without any specific knowledge on the dataset. As we optimize the transformation with respect to the downstream biological objective, this parametric framework reweighs the contribution of each hidden factor and makes the biological signals visible. We demonstrate that DataRemix can outperform normalization methods which make explicit use of dataset specific technical factors. Also we show that DataRemix can be efficiently optimized via Thompson Sampling approach, which makes it feasible for computationally expensive objectives such as eQTL analysis. Finally we reanalyze the Depression Gene Networks (DGN) dataset, and we highlight new *trans*-eQTL networks which were not reported in the initial study.

1 Genome-wide gene expression studies have become  
2 a staple of large scale systems biology and clinical  
3 projects. However, while gene expression is the most  
4 mature high-throughput technology, technical chal-  
5 lenges remain. Raw gene expression values must be  
6 normalized for any technical and nuisance biological  
7 variation and the normalization strategy can have dra-  
8 matic effects on the results of downstream analysis.  
9 This is especially true in cases where the sought-  
10 after gene expression effects are likely to be small  
11 in magnitude, such as expression quantitative trait  
12 loci (eQTLs). Increasingly sophisticated normalization  
13 methods have been proposed and many are computa-  
14 tional intensive and/or can have multiple free param-  
15 eters that must be optimized (Leek & Storey 2007;  
16 Stegle *et al.* 2010; Listgarten *et al.* 2010; Kang *et al.*  
17 2008; Mostafavi *et al.* 2013). Moreover, it is not un-  
18 common for one dataset to yield multiple normalized  
19 versions that maximize performance in a particular  
20 setting (such as the discovery of *cis*- and *trans*-eQTLs  
21 Battle *et al.* 2014), highlighting the complexity of the  
22 normalization problem.

23 Singular value decomposition (SVD) is one of the  
24 most widely used gene expression analysis tools (Al-  
25 ter *et al.* 2000, 2003) that can also be used for data  
26 normalization. Using the SVD we can simply remove  
27 the first few principle components that are presumed  
28 to represent technical factors such as batch-effects or  
29 other nuisance variation. In some cases this dramati-  
30 cally improves downstream performance, for example

in the case of eQTL analysis (Mostafavi *et al.* 2013).  
The drawback of this method is that the exact number  
of components to remove must be determined empiri-  
cally and some meaningful biological signals may be  
lost in the process.

More sophisticated approaches attempt to partition  
data structure into true biological and nuisance varia-  
tion and remove only the latter (Leek & Storey 2007;  
Stegle *et al.* 2010; Listgarten *et al.* 2010; Kang *et al.*  
2008; Mostafavi *et al.* 2013). These can improve on  
the naive SVD-based normalization but require addi-  
tional input such as technical covariates, or the study  
design. The success of these methods ultimately de-  
pends on the availability and quality of such meta data  
and some methods still rely on parameter optimization  
to maximize performance. These widely used normal-  
ization approaches all have a common theme that the  
rely in part on the intrinsic data structure. One key  
property that contributes to the success of these ap-  
proaches is that for many biological questions of inter-  
est nuisance variation (of technical or biological origin)  
is larger in magnitude than true biological variation.  
Our proposed method, DataRemix, explicitly formal-  
izes this view of the data normalization problem.

In this work we demonstrate that biological util-  
ity of gene expression datasets can be dramatically  
improved with a simple three-parameter transforma-  
tion, DataRemix. Our method does not require any  
dataset specific knowledge but rather optimizes the  
transformation with respect to some independent *ob-*  
*jective* of data quality, such as the quality of the gene-  
correlation network or the number of *trans*-eQTL dis-  
coveries. Because our method requires only the gene  
expression data and biological validity objective, it can

\*Correspondence: [mchikina@pitt.edu](mailto:mchikina@pitt.edu)

<sup>2</sup>Department of Computational and Systems Biology, School of Medicine,  
University of Pittsburgh,

Full list of author information is available at the end of the article

65 be applied to any publicly available dataset. We focus  
 66 our study on gene expression data for which methods  
 67 for quantifying biological validity are well established,  
 68 but our approach can be readily applied to any high-  
 69 throughput molecular data for which similar quality  
 70 metrics can be defined. We show that this strategy can  
 71 outperform methods that make explicit use of dataset  
 72 specific factors, and can further improve datasets that  
 73 have been extensively normalized via an optimized, pa-  
 74 rameter rich model. We also show how the optimal  
 75 parameters of DataRemix can be found efficiently by  
 76 Thompson Sampling with a dual learning setup, mak-  
 77 ing the approach feasible for computationally expen-  
 78 sive objectives such as eQTL analysis.

## 79 Result

### 80 The DataRemix framework

We formulate DataRemix as a simple parametrized  
 version of SVD which can be directly optimized to  
 improve the biological utility of gene expression data.  
 Given a gene-by-sample matrix  $X$ , SVD decomposi-  
 tion can be thought of as a solution to the low-rank  
 matrix approximations problem defined as:

$$\min_{U_k, \Sigma_k, V_k} \|X - U_k \Sigma_k V_k^T\|_F^2 \quad (1)$$

where  $U$  and  $V$  are unitary matrices. With the SVD  
 decomposition  $U \Sigma V^T$ , the product of  $k$ -truncated ma-  
 trices  $U_k \Sigma_k V_k^T$  gives the rank- $k$  approximation of  $X$ .  
 We introduce two additional parameters  $p$  and  $\mu$  to  
 define a new reconstruction:

$$\text{DataRemix}_{\{k,p,\mu\}}(X) = U_k \Sigma_k^p V_k^T + \mu(X - U_k \Sigma_k V_k^T) \quad (2)$$

81 Here,  $k$  is the number of principle components of SVD  
 82 and  $p \in [-1, 1]$  is a real number which alters the scaling  
 83 of each eigenvalue. For  $p = 1$ , this approach reduces  
 84 to the original SVD-based reconstruction. For  $p = 0$   
 85 the transformation gives the frequently used whiten-  
 86 ing operation (Friedman 1987). As depicted in Figure  
 87 1, generally, different choices of  $p$  reweigh the con-  
 88 tribution of each variance component, possibly mak-  
 89 ing some low-variance biological signals visible while  
 90 down-weighting technical and other systematic noise.  
 91 The parameter  $\mu$  is a non-negative weight that adds  
 92 the residual back to the reconstruction in order to  
 93 make the transformation *lossless*.

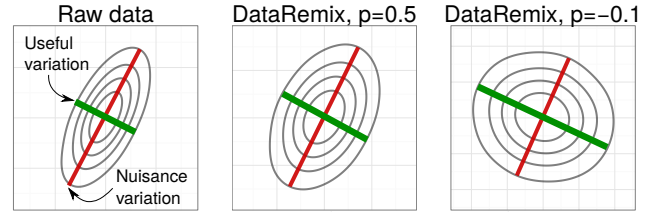


Figure 1: Visual representation of DataRemix transformation. We simulate a 2-dimensional dataset where the nuisance variation contributes more variance than true biological variation. Different power parameters  $p$  reweigh the contributions of the two variance axes, making the true biological variation more “visible”.

Intuitively, we expect this approach to succeed be-  
 cause sophisticated normalization methods that use  
 both data structure and some external variables, such  
 as technical covariates, can be thought of as implicit  
 regularizations on the naive SVD-based normalization  
 (which simply removes the first  $k$  components), and  
 this formulation simply makes this explicit.

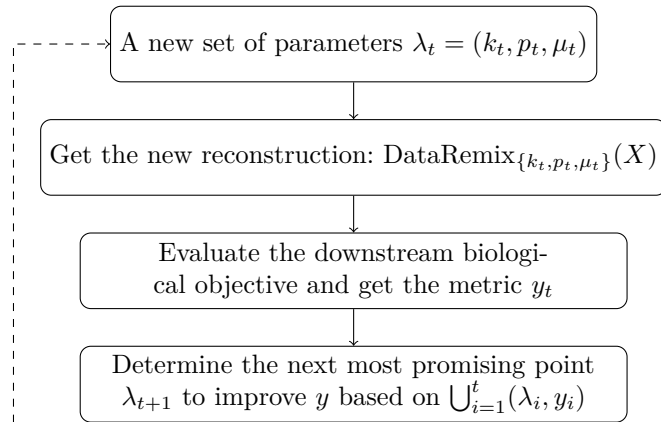


Figure 2: The workflow of DataRemix.

The general workflow of DataRemix is shown in Fig-  
 ure 2. The downstream biological objective depends  
 on your study. For example, if you focus on *trans*-eQTL  
 analysis, the biological objective will be to increase the  
 number of *trans*-eQTLs detected from the DataRemix-  
 normalized gene expression profile and the metric  $y$   
 will be the number of *trans*-eQTLs deemed significant.  
 The parameter optimization step which determines the  
 next point to check is detailed in the next section.

### Parameter Optimization

The parameters  $\lambda = (k, p, \mu)$  need to be optimized  
 with respect to a particular biological objective. Grid  
 search and random search (Bergstra & Bengio 2012)  
 are among the most popular strategies, but these  
 methods have low efficiency. Most of the search steps  
 are wasted and the optimality of parameters is highly  
 constrained by the step size and available computing  
 power. In order to utilize the search history and keep

121 a good balance between exploration and exploitation,  
 122 we can formulate parameter search as a dual learning  
 123 task.

124 We define a general performance measure  $y =$   
 125  $L(\lambda, \mathcal{D})$ , with  $\lambda$  representing the parameter tuple  
 126  $(k, p, \mu)$ ,  $\mathcal{D}$  as the data,  $L$  as the evaluating process and  
 127  $y$  as the biological objective. Ideally we can determine  
 128 the optimal point  $\operatorname{argmax}_{\lambda} L$  easily by gradient descent  
 129 based method, but usually  $L$  is derivative-free and it  
 130 is time intensive. Thus we introduce a surrogate model  
 131  $f(\lambda)$  which can directly predict  $L(\lambda, \mathcal{D})$  only given  $\lambda$ .  
 132 There are two conditions on  $f$ :  $\operatorname{argmax}_{\lambda} f$  should be  
 133 easy to solve and  $f$  should have enough capacity.

With these two properties, we can sequentially  
 update  $f$  with  $(\lambda_t, y_t)$  and propose to evaluate  $L$   
 at  $\lambda_{t+1} = \operatorname{argmax}_{\lambda} f$  in the next step. By gradu-  
 ally updating  $f$  with newly evaluated samples  $(\lambda, y)$ ,  
 $\operatorname{argmax}_{\lambda} f$  approaches the true underlying optimal  
 $\operatorname{argmax}_{\lambda} L$  as  $f$  can gradually fit to the underlying  
 mapping function  $L$ . This provides a more efficient  
 approach to explore the parameter space by exploit-  
 ing the search history. In this work, we model  $f$  as  
 a sample from a Gaussian Process with mean 0 and  
 kernel  $k(\lambda, \lambda')$ , where  $\lambda = (k, p, \mu)^T$ . It is well known  
 that the form of the kernel has considerable effect on  
 performance. After experimentation we settled on the  
 exponential kernel as the most suited for our applica-  
 tion. The exponential kernel is defined as below (note  
 the difference from the squared-exponential or RBF  
 kernel).

$$k(\lambda, \lambda') = \exp\left(-\frac{\|\lambda - \lambda'\|_2}{2}\right) \quad (3)$$

134 We observe  $y_t = f(\lambda_t) + \epsilon_t$ , where  $\epsilon_t \sim N(0, \sigma^2)$ . For  
 135 Bayesian optimization, one approach for picking the  
 136 next point to sample is to utilize acquisition functions  
 137 (Snoek *et al.* 2012) which are defined such that high  
 138 acquisitions correspond to potentially improved per-  
 139 formance. An alternative approach is the Thompson  
 140 Sampling approach (Basu & Ghosh 2017; Agrawal &  
 141 Goyal 2013; Hernández-Lobato *et al.* 2014). After we  
 142 update the the posterior distribution  $P(f|\lambda_{1:t}, y_{1:t})$ , we  
 143 draw one *sample*  $f$  from this posterior distribution as  
 144 the optimization target to infer  $\lambda_{t+1}$ . Theoretically it  
 145 is guaranteed that  $\lambda_t$  converges to the optimal point  
 146 gradually (Agrawal & Goyal 2013). With this theoret-  
 147 ical guarantee, we focus on Thompson Sampling ap-  
 148 proach to optimize parameters for DataRemix.

#### 149 Estimation of Hyper-Parameters

150 First we rely on the maximum likelihood estimation  
 151 (MLE) to infer the variance of noise  $\sigma^2$  (Rasmussen  
 152 2004). Given the marginal likelihood defined by (4), it

is easy to use any gradient descent method to deter-  
 mine the optimal  $\sigma^2$

$$\log p(\vec{y}|\vec{\lambda}) = -\frac{1}{2}\vec{y}^T(K + \sigma^2 I)^{-1}\vec{y} - \frac{1}{2}\log|K + \sigma^2 I| - \frac{t}{2}\log 2\pi \quad (4)$$

where  $\vec{y} = y_{1:t} = (y_1, \dots, y_t)^T$ ,  $\vec{\lambda} = \lambda_{1:t} = (\lambda_1, \dots, \lambda_t)^T$   
 and  $K$  is the covariance matrix with each entry  
 $K_{ij} = k(\lambda_i, \lambda_j)$ .

#### 158 Sampling from the Posterior Distribution

159 Since Gaussian Process can be viewed as Bayesian  
 160 linear regression with infinitely many basis functions  
 161  $\phi_0(\lambda), \phi_1(\lambda), \dots$  given a certain kernel (Rasmussen  
 162 2004), in order to construct an analytic formulation  
 163 for the sample  $f$ , first we need to construct a certain  
 164 set of basis functions  $\Phi(\lambda) = (\phi_0(\lambda), \phi_1(\lambda), \dots)$ , which  
 165 is also defined as feature map of the given kernel. Then  
 166 we can write the kernel  $k(\lambda, \lambda')$  as the inner product  
 167  $\Phi(\lambda)^T \Phi(\lambda')$ .

Mercer's theorem guarantees that we can express the  
 kernels in terms of eigenvalues and eigenfunctions, but  
 unfortunately there is no analytic solution given the  
 exponential kernel we used. Instead we make use of the  
 random Fourier features to construct an approximate  
 feature map (Rahimi & Recht 2008). First we compute  
 the Fourier transform  $p$  of the kernel (see Supplemental  
 Note for derivation).

$$p(\vec{\omega}) = \frac{1}{(2\pi)^3} \int \exp(-i\vec{\omega}^T \vec{\Delta}) \exp\left(-\frac{\|\vec{\Delta}\|_2}{2}\right) d\vec{\Delta} \quad (5)$$

$$= \frac{8}{\pi^2(4\|\vec{\omega}\|_2^2 + 1)^2}$$

where  $\vec{\omega} = (\omega_1, \omega_2, \omega_3)^T$  and  $\vec{\Delta} = \lambda - \lambda'$ . Then we  
 draw  $m_t$  iid samples  $\omega_1, \dots, \omega_{m_t} \in \mathbb{R}^3$  by rejection  
 sampling with  $p(\omega)$  as the probability distribution.  
 Also we draw  $m_t$  iid samples  $b_1, \dots, b_{m_t} \in \mathbb{R}$  from  
 the uniform distribution on  $[0, 2\pi]$ . Then the feature  
 map is defined by the following equation.

$$\Phi(\lambda) = \sqrt{\frac{2}{m_t}} [\cos(\omega_1^T \lambda + b_1), \dots, \cos(\omega_{m_t}^T \lambda + b_{m_t})]^T \quad (6)$$

where the dimension  $m_t$  can be chosen to achieve the  
 desired level of accuracy with respect to the difference  
 between true kernel values  $k(\lambda, \lambda')$  and the approxi-  
 mation  $\Phi(\lambda)^T \Phi(\lambda')$ .

172 *Thompson Sampling*

Any sample  $f$  from the Gaussian Process can be defined by  $f(\lambda) = \Phi(\lambda)^T \theta$ , where  $\theta \sim N(0, I)$  and  $\Phi(\lambda)^T$  is defined by (6). In order to draw a posterior sample  $f$ , we just need to draw a random sample  $\theta$  from the posterior distribution  $P(\theta | \vec{\lambda}, \vec{y})$ .

$$P(\theta | \vec{\lambda}, \vec{y}) \propto P(\vec{y} | \vec{\lambda}, \theta) P(\theta) \tag{7}$$

$$\propto N(A^{-1} \Phi(\vec{\lambda}) \vec{y}, \sigma^2 A^{-1})$$

173 where  $A = \Phi(\vec{\lambda}) \Phi(\vec{\lambda})^T + \sigma^2 I$  and  $\Phi(\vec{\lambda}) = (\Phi(\lambda_1) \cdots \Phi(\lambda_t))$   
 174 (see Supplemental Note for more details). The overall  
 175 algorithm is summarized as the following pseudo code.

**Algorithm 1** Thompson Sampling for Searching  $\lambda$

Extra Parameters

$t_{max}$ : the maximum number of iteration steps

$\xi$ : a pre-defined probability which ensures the search doesn't get stuck in a local optimum

1. Get a short sequence  $\mathcal{D}_1 = (\lambda, y)$  as seeds by random search.
2. Draw  $m_t$  iid samples  $\omega_1, \dots, \omega_{m_t} \in \mathbb{R}^3$  and  $m_t$  iid samples  $b_1, \dots, b_{m_t} \in \mathbb{R}$  according to (5)
3. Iterate from  $t = 1$  until  $\lambda$  converges or it reaches  $t_{max}$ 
  - (1) At step  $t$ , estimate the hyper-parameter  $\sigma^2$  given  $\mathcal{D}_t$  according to (4)
  - (2) Draw a sample  $f$  given  $\mathcal{D}_t$  according to (7) with feature map determined by (6)
  - (3)  $\lambda_{t+1} = \begin{cases} \operatorname{argmax}_{\lambda} f(\lambda) & \text{w.p. } 1 - \xi \\ \text{random search} & \text{w.p. } \xi \end{cases}$
  - (4) Evaluate  $y_{t+1}$  given  $\lambda_{t+1}$
  - (5)  $\mathcal{D}_{t+1} = \mathcal{D}_t \cup (\lambda_{t+1}, y_{t+1})$

176 Quality of the correlation network derived from the  
 177 GTex gene expression study.

178 The GTex datasets (Lonsdale *et al.*, 2013) is comprised  
 179 of human samples from diverse tissues, many of which  
 180 were obtained post-mortem and there are many technical  
 181 factors which have considerable effects on the gene  
 182 expression measurements. On the other hand this rich  
 183 dataset provides an unprecedented multi-tissue map of  
 184 gene regulatory networks and has been extensively analyzed  
 185 in this context. It is natural to assume that a  
 186 dataset that is better at recovering known pathways is  
 187 likely to yield more credible novel predictions. Thus,  
 188 we use DataRemix to optimize the known pathway recovery  
 189 task as a function of the correlation network computed on a  
 190 Remixed dataset.

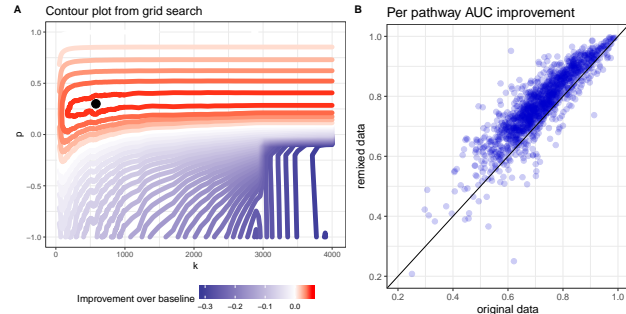


Figure 3: **A** The improvement in performance of DataRemix transform of the pathway prediction task visualized as a function of  $k$  and  $p$  parameters ( $\mu$  is fixed at 0.01). Performance is measured as the mean AUC across all pathways in the “canonical” mSigDB dataset and the red contours indicate improvement over the performance on untransformed data. **B** Per-pathway performance improvement for the optimal DataRemix transformation.

Specifically we start with a quantile normalized TPM data that has not been corrected for technical factors or tissue of origin. We formally define the objective as the average AUC across “canonical” mSigDB pathways (which include KEGG, Reactome and PID) (Subramanian *et al.*, 2005) using guilt-by-association. Specifically, the genes are ranked by their average Pearson correlation to other genes in the pathway (excluding the gene when the gene itself is a pathway member). Figure 3A depicts the results of grid search for the parameters  $k, p$  (with  $\mu$  fixed at 0.01) and the contour plot shows a clear region of increased performance. Using the optimal transformation found by grid search, we plot per-pathway AUC improvement in Figure 3B and find that the AUC is substantially increased for almost every pathway.

eQTL discovery in the DGN dataset.

We also consider the task of discovering *cis*- and *trans*-eQTLs on the Depression Gene Networks (DGN) dataset (Battle *et al.*, 2014). In the original analysis this dataset was normalized using the Hidden Covariates with Prior (HCP) (Mostafavi *et al.*, 2013) with four free parameters that were separately optimized for *cis*- and *trans*-eQTLs. The rationale behind separate *cis* and *trans* optimized normalization can be understood in terms of which variance components represent true biological vs. nuisance variation in the two contexts. Specifically, *cis*-eQTLs represent *direct* effects of genetic variation on the expression of a single gene. On the other hand, *trans*-eQTLs represent network level, *indirect* effects that are mediated by a regulator. Thus, *trans*-eQTLs are reflected in systematic variation in the data which becomes a nuisance factor when only direct effects are of interest. It thus follows that the data should be more aggressively normalized for *cis*-eQTL discovery. The original analysis of this dataset optimized the HCP parameters separately for the *cis*

191

192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207

208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228



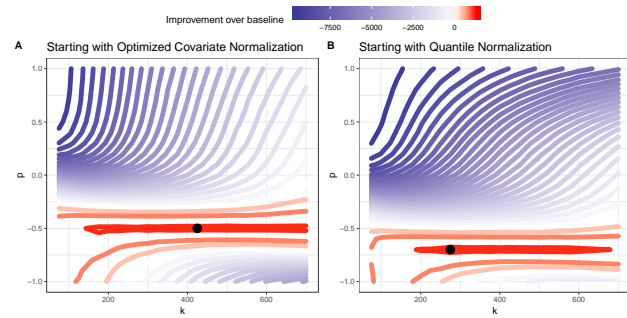
229 and *trans* tasks yielding two different datasets that we  
 230 refer to as  $D_{cis-optim}$  and  $D_{trans-optim}$ .

231 The HCP model takes various technical covariates as  
 232 input, and of the covariates used in the original study  
 233 20 cannot be inferred from the gene-level counts. In  
 234 order to investigate how much improvement can be  
 235 achieved via DataRemix in the absence of access to  
 236 these covariates we also consider a “naively” normal-  
 237 ized dataset, quantile normalization of log-transformed  
 238 counts, or  $D_{QN}$ .

239 *cis-eQTLs*.

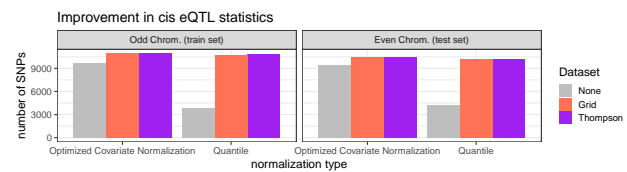
240 In this task we focus on optimizing the discovery of  
 241 *cis*-eQTLs. We define *cis*-eQTLs as a SNP-gene in-  
 242 teraction where the SNP is located within 50kb of the  
 243 gene’s transcription start site. The interaction is quan-  
 244 tified with Spearman rank correlation and deemed sig-  
 245 nificant at 10% FDR (Benjamini-Hochberg correction  
 246 for the total number of tests).

247 We perform our analysis in a cross-validation frame-  
 248 work, whereby we can optimize DataRemix param-  
 249 eters (using grid search or Thompson Sampling) using  
 250 SNPs on the odd chromosomes only and then evaluate  
 251 the parameters on the held-out even chromosome set.  
 252 We visualize the effect of varying the  $k$  and  $p$  param-  
 253 eters on the performance of the DataRemix transform  
 254 in Figure 4. Red regions indicate improvement over the  
 255 number of *cis*-eQTLs discovered with the  $D_{cis-optim}$   
 256 dataset. We find that both versions of the dataset can  
 257 be improved via the DataRemix transform to a simi-  
 258 lar degree. We also find that on this task the optimal  $p$   
 259 parameter is negative and the result is relatively insen-  
 260 sitive to the choice of  $k$ . The last observation can be  
 261 interpreted when we consider the interaction between  
 262  $p$  and  $\mu$  (the multiplier for the residual part including  
 263  $k+1$  through  $\max(k)$  components). If we wish to bring  
 264 forward small-variance components, as is the case with  
 265 *cis*-eQTL discovery, we would like the diagonal values  
 266 of  $\mu \Sigma_{k+1:\text{rank of } X}$ , representing the contribution of the  
 267 later components, to be in the same range or larger  
 268 than  $\max(\Sigma_{1:k}^p)$  which is the largest contribution of  
 269 the high variance components. This can be achieved  
 270 by picking different values of  $k$ .



271 Figure 4: Contour plot representing the effects of the  $k$  and  $p$  pa-  
 rameters on the performance of DataRemix on *cis*-eQTL discov-  
 ery on 50,000 randomly selected SNPs on odd chromosomes (train-  
 ing set). Red contours represent parameter combinations that in-  
 crease the number *cis*-eQTLs beyond what can be achieved using  
 the  $D_{cis-optim}$  dataset. Panel A shows the results starting with  
 $D_{cis-optim}$  while  $D_{QN}$  is used for panel B. Improvement can be  
 achieved starting with either dataset. We note that the optimal  $p$   
 parameter is negative (though slightly different) for both datasets.

272 The final results for both the train and test set  
 273 are depicted in Figure 5. We find that the optimal  
 274 parameters are indeed generalizable as we achieve a  
 275 similar level of improvement on the train and test  
 276 datasets. Importantly, we find that while the quantile-  
 277 normalized dataset  $D_{QN}$  performs considerably worse  
 278 than  $D_{cis-optim}$  the two datasets achieve comparable  
 279 performance after applying DataRemix. Moreover, the  
 280 final performance of the Remixed  $D_{QN}$  dataset is an  
 281 improvement of the baseline  $D_{cis-optim}$  demonstrat-  
 282 ing the near optimal normalization is possible with-  
 283 out access to technical covariates. We do note, on this  
 284 task, the final performance of the Remixed  $D_{cis-optim}$   
 285 is slightly better than that of  $D_{QN}$  and thus it is still  
 286 advisable to include such covariates in the normaliza-  
 287 tion pipeline if they are available.



288 Figure 5: Final results from DataRemix parameter search using a  
 cross-validation framework. Optimal parameters are determined us-  
 ing the odd chromosome SNPs only and then tested on the even  
 chromosome SNPs. We find that the DataRemix transform does  
 not overfit the objective as the degree of improvement is similar  
 across the test and train SNP sets (note: the starting value of the  
 baseline (DataRemix=“None”) datasets differ between the test and  
 train SNP set). Moreover, we find that Thompson Sampling is able  
 to match grid search results using only 100 evaluations.

289 *trans*-eQTLs.

290 In our third task, we optimize the discovery of *trans*-  
 291 eQTLs in the same DGN dataset. Ideally, *trans*-eQTLs  
 292 represent network-level effects and thus give some in-  
 293 sight about the regulatory structure of gene expres-  
 294 sion. However, in practice *trans*-eQTLs are simply defined  
 295 as SNP-gene associations where the SNP and

the gene are located on different chromosomes. While this is a useful heuristic definition, it doesn't guarantee that the association is mediated at the network level. One possible source of bias is mis-mapped RNAseq reads which contaminate the quantification of the apparently *trans*-associated gene with reads from a homologous locus that has *cis* association. Even in the absence of technical artifacts, direct interchromosomal interactions have been observed (see Williams *et al.* 2010 for a comprehensive review). In order to focus on potential indirect effects, we apply an additional filter to *trans*-eQTL discovery. Specifically we require SNPs involved in a *trans* effect to be associated with more than one gene at a FDR of 20% (Benjamini-Hochberg correction for the total number of test (approximately  $8 \times 10^9$ ). We term these SNPs *trans*-SNPs<sup>+</sup>. In comparison with same chromosome *cis*-eQTLs, inter-chromosome *trans*-eQTLs are rare and *trans*-SNPs<sup>+</sup> (as defined above) are more rare still. In fact, using the odd chromosome SNPs subsampled at 20%, we find only 88 such SNPs using  $D_{trans-optimal}$  dataset and this is the default value we wish to improve.

As is the case with *cis*-eQTLs, we investigate the  $k, p$  performance surface of the DataRemix transform at the grid-search optimal  $\mu = 0.01$ . Given that the relevant variance components that would maximize the *trans*-eQTL objective are different, it is not surprising that we find that the performance surface differs as well. In particular, we find that the optimum value of  $p$  is positive but close to 0 and thus the first  $k$  variance components are weighted equally with a weight close to 1. Consequently, at  $\mu = 0.01$  and  $p \approx 0$  the contribution of the first  $k$  components is considerably larger than that of the remaining ones and we find that the performance is more sensitive to the exact value of  $k$ .

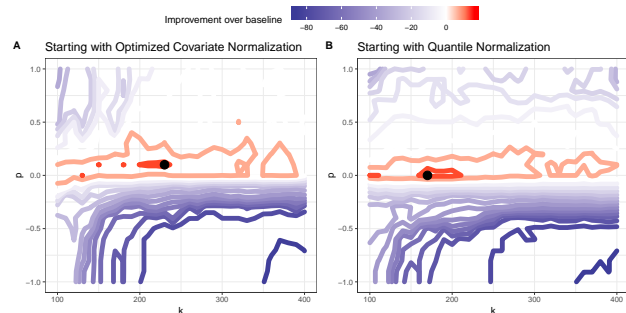


Figure 6: Contour plot representing the effects of the  $k$  and  $p$  parameters on the performance of DataRemix on *trans*-eQTL discovery on 50,000 randomly selected SNPs on odd chromosomes (training set). Red contours represent parameter combinations that increase the number of *trans*-eQTLs beyond what can be achieved using the  $D_{trans-optimal}$  dataset. Panel A shows the results starting with  $D_{trans-optimal}$  while  $D_{QN}$  is used for panel B. Improvement can be achieved starting with either datasets. We note that the performance is more sensitive to the choice of  $k$ .

Despite the difference in the performance landscape, we find that the DataRemix transform behaves similarly on this objective. Specifically, either starting dataset can be improved to similar final performance, though the optimal parameters are slightly different. As is the case with the *cis*-eQTL objective, the cross-validation procedure gives consistent results and no overfitting is observed for either grid search or Thompson Sampling (Figure 7).

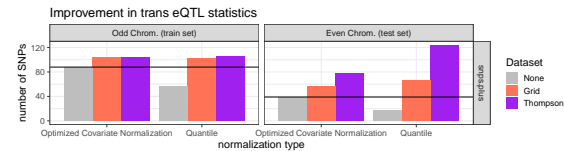
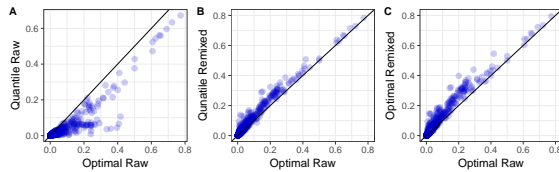


Figure 7: Final values for the eQTL statistics obtained from two versions of datasets. Here we make a comparison between quantile normalized  $D_{QN}$  and HCP normalized  $D_{trans-optimal}$  with parameters optimized for *trans*-eQTL discovery. We find DataRemix is able to improve upon either of starting datasets and the improvement on both the train and test dataset are comparable which indicates that overfitting is not a problem

Since *trans*-eQTLs are likely to reflect pathway level effects, we expect that a dataset that is optimally transformed for *trans*-eQTL discovery should also produce better correlation networks. We thus investigate if optimal DataRemix transform is transferable between tasks by checking if Remixed dataset optimized with respect to *trans*-eQTL discovery also improves the network quality criterion. Similar to our analysis of the GTex datasets, we use the correlation network to perform guilt-by-association pathway predictions and evaluate the results over 1,330 MSigDB canonical pathways. Figure 8 shows scatter plots of per-pathway AUPR (area under precision-recall curve) for several comparisons with respect to the baseline  $D_{trans-optimal}$  dataset. In the first panel we contrast the performance to  $D_{QN}$  and we observe that  $D_{trans-optimal}$  brings a considerable improvement over the quantile normalized dataset. In the second panel we contrast  $D_{trans-optimal}$  with the Remixed version of  $D_{QN}$  (optimized for *trans*-eQTL discovery with Thompson Sampling). We find that the pattern becomes opposite and the Remixed  $D_{QN}$  dataset performs consistently better than  $D_{trans-optimal}$ . The final panel shows the results of Remixing  $D_{trans-optimal}$  itself which also improves the performance. Overall, we find that DataRemix improves multiple criteria of biological validity as optimizing for the *trans*-eQTL objective also results in improved correlation networks. Interestingly, we find that while the Remixed  $D_{trans-optimal}$  is no better than Remixed  $D_{QN}$  on *trans*-eQTL discovery, it performs slightly better on the pathway prediction task. Taking the two objectives into account, we conclude that starting with a properly covariate-normalized dataset

376 is superior overall, which is also the our finding regard-  
 377 ing the *cis*-eQTL objective.



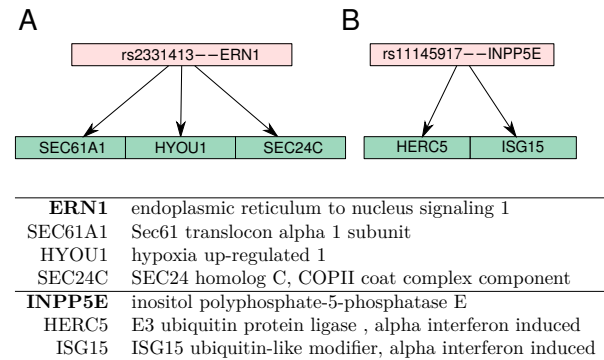
378 Figure 8: DataRemix-transformed datasets improve the pathway  
 prediction objective which is not explicitly optimized. Each plot  
 is a per-pathway AUPR (area under precision-recall curve) from  
 various datasets (y-axis) contrasted with the results from the opti-  
 mal covariate-normalized dataset  $D_{trans-optimal}$ , which serves as the  
 baseline (x-axis). Panel A shows the contrast between  $D_{trans-optimal}$   
 and  $D_{QN}$ . The performance of  $D_{trans-optimal}$  is considerably better.  
 Panel B shows the results of the Remixed  $D_{QN}$  datasets (optimized  
 for *trans*-eQTL discovery with Thompson Sampling). Even though  
 $D_{QN}$  starts out as considerably worse, the Remixed version is able  
 to outperform  $D_{trans-optimal}$ . Panel C shows the results of Remixed  
 $D_{trans-optimal}$ . We choose to use AUPR instead of AUC because we  
 find that Remixed version matches but doesn't further improve the  
 AUC performance of  $D_{trans-optimal}$

379 A major finding of our study is that for the eQTL and  
 380 pathway prediction tasks, the starting point of nor-  
 381 malizing DGN datasets appears to matter relatively  
 382 little. Even though the quantile-normalized dataset  
 383 performs considerably worse in the beginning, after  
 384 Remixing its performance matches that of the opti-  
 385 mal covariate-normalized datasets. Of course, if covari-  
 386 ates are available, it is preferable to use them and in  
 387 the case of DGN, slightly further improvement can be  
 388 achieved. However our results indicate that in some  
 389 cases datasets *can* be effectively normalized even in the  
 390 absence of meta-data about quality control or batch  
 391 variables which is an important consideration for many  
 392 legacy datasets where such information is not avail-  
 393 able.

### 394 Novel Biological Findings

395 At the optimal DateRemix parameters for  $D_{QN}$ , we  
 396 find an additional 24 loci that have significant associ-  
 397 ations with more than one gene and are not in link-  
 398 age disequilibrium with those significant hits in the  
 399  $D_{trans-optimal}$ . We highlight two examples of new regu-  
 400 latory modules recovered via DataRemix that ap-  
 401 pear to be biologically credible based on the known  
 402 functions of the genes involved. One of the newly sig-  
 403 nificant interactions involves the SNP rs2331413 lo-  
 404 cated in proximity of the ERN1 gene, which func-  
 405 tions as a sensor of unfolded protein in the endoplas-  
 406 mic reticulum and triggers an intracellular signalling  
 407 pathway termed the unfolded protein response. Three  
 408 downstream genes associated with rs2331413 are like-  
 409 wise endoplasmic reticulum proteins. The ERN1 lo-  
 410 cus has been associated with several phenotypes in  
 411 GWAS studies, most notably drug induced hepatotox-  
 412 icity (Petros *et al.* 2017).

We also find an SNP rs11145917 located near  
 INPP5E gene which is associated with two genes in  
 the alpha interferon response. Even though only two  
 genes show genome-wide significance, several other  
 canonical members of the alpha interferon response  
 are just slightly short of the significance threshold sug-  
 gesting that the locus affects the upstream signaling  
 components. The INPP5E locus has been implicated  
 in a variety of autoimmune diseases as well as blood  
 immune-cell composition phenotype (de Lange *et al.*  
 2017; Astle *et al.* 2016), though to our knowledge no  
 mechanism has been proposed. Our analysis suggests  
 that INPP5E may affect baseline activity of the alpha  
 interferon pathway, which is a testable prediction with  
 potential clinical importance.



428 Figure 9: Clusters of *trans*-eQTLs detected by DataRemix that were  
 not significant in the original dataset. Panel A. Both the *cis* and  
*trans* genes are involved in ER biology and specifically unfolded  
 protein response. Panel B. Both of the *trans* genes are canonical  
 targets of alpha interferon. The upstream *cis* gene, INPP5E, is a  
 signaling molecule that mediates cell responses to various stimula-  
 tion and its locus has been implicated in a variety of autoimmune  
 diseases as well as blood immune-cell composition phenotypes.

### 429 Thompson Sampling Performance

430 We find that Thompson Sampling matches the best  
 431 grid-search performance in under 100 steps giving a 40-  
 432 fold reduction in the number of evaluations. We also  
 433 note that it is possible for the Thompson sampling  
 434 to surpass the grid-search results since the parameter  
 435 combinations are not constrained by the choice of grid.

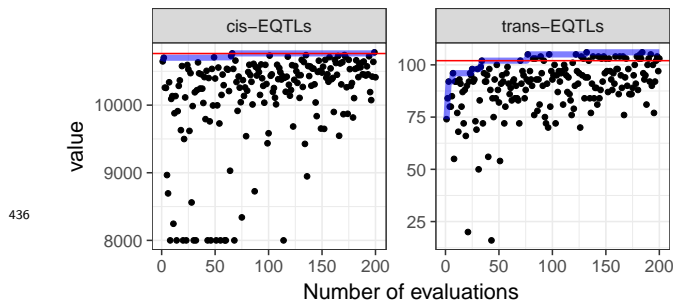


Figure 10: Objective evaluations as a function of iteration number for the *trans*-eQTL and *cis*-eQTL objectives using the quantile normalized  $D_{QN}$  dataset. Red lines indicate the maximum value that was obtained by grid-search and blue lines indicate the cumulative maximum of Thompson Sampling.

## Discussion

We have proposed DataRemix, a new optimizable transformation for gene expression data. The transformation is able to improve the biological validity of gene expression representations and can be used for effective normalization in the absence of any knowledge of technical covariates. One limitation of the DataRemix approach is that it works best on data that is well approximated by a single Gaussian. However, it is relatively straightforward to adapt the approach to matrix decompositions different from SVD that are more suitable for non-Gaussian data, such as independent component analysis. We also note that it is possible to introduce additional parameters that specify more complex weighting schemes. However, as the number of parameters is increased, there is a potential for over-optimization of a specific objective above others. We emphasize that in our simple parametrization, we observe that multiple metrics of biological validity improve when only one is explicitly optimized. Specifically we find that optimizing for *trans*-eQTL discovery also improves the correlation network as measured by guilt-by-association pathway prediction. This property is less likely to be preserved as the number of parameters is increased.

## Methods

### GTex Dataset

We downloaded the complete gene-level TPM data (RNASeqCV1.1.8) from the GTex consortium (Lonsdale *et al.*, 2013). These data were quantile normalized.

### DGN Dataset

Depression Gene Networks (DGN) dataset contains whole-blood RNA-seq and genotype data from 922 individuals. The genotype data was filtered for  $MAF > 0.05$ . The genomic coordinate of each SNP was taken from the Ensembl Variation database (version 90, hg19/GRCh37). SNP identifiers that were not present

in that release were excluded. After filtering there were 649,875 autosomal single nucleotide polymorphisms (SNPs). Data is available upon application through NIMH Center for Collaborative Genomic Studies on Mental Disorders. For gene expression we used the gene-level quantified dataset. The dataset comes already filtered for expressed genes and was further filtered for gene symbols that were not present in Ensembl 90 leaving 13,708 genes. The dataset comes in two covariate normalized versions with normalization parameters optimized for *cis*- and *trans*-eQTL discovery separately. To create the naive-normalized dataset, we applied a log transformation,  $\log(x + 1)$ , to the raw counts and quantile normalized the results.

### eQTL mapping

eQTL association mapping was quantified with Spearman rank correlation. For *cis*-eQTLs, testing was limited to SNPs which locate within 50kb of any of the gene's transcription start sites (Ensembl, version 90). *cis*-eQTL is deemed significant at 10% FDR with Benjamini-Hochberg correction for the total number of tests. For *trans*-eQTLs, the significance cutoff is 20% FDR with Benjamini-Hochberg correction for the total number of tests. Since the Benjamini-Hochberg FDR is a function of the entire p-value distribution in order to ensure consistency comparisons, the rejection level was set once based on the p-value that corresponded to 10% or 20% FDR in the original *cis*-optimized  $D_{cis-optimal}$  and *trans*-optimized  $D_{trans-optimal}$  dataset respectively. To reduce the computational cost of grid evaluations, all the optimization computations were performed on a set of 100,000 subsampled SNPs.

### Correlation network evaluation

We evaluated the quality of the correlation network derived from a particular dataset using guilt-by-association pathway prediction. Specifically, the genes were ranked by their average Pearson correlation to other genes in the pathway (excluding the gene when the gene itself is a pathway member). The resulting ranking was evaluated for performance using AUC or AUPR metric. For pathway ground-truth we used the "canonical" pathways dataset from MSigDB, comprising 1,330 pathways (Subramanian *et al.*, 2005).

### Software Access

DataRemix is an R package which is freely available at GitHub (<https://github.com/wgmao/DataRemix>).

### Competing interests

The authors declare that they have no competing interests.

### Author's contributions

Text for this section ...



524 **Acknowledgements**

525 Text for this section . . .

526 **Author details**

527 <sup>1</sup>Joint Carnegie Mellon-University of Pittsburgh Ph.D. Program in  
528 Computational Biology,. <sup>2</sup>Department of Computational and Systems  
529 Biology, School of Medicine, University of Pittsburgh,.

530 **References**

531 Leek JT, Storey JD. 2007. Capturing heterogeneity in gene expression  
532 studies by surrogate variable analysis. *PLoS Genet*, **3**: 1724–1735.  
533 Stegle O, Parts L, Durbin R, Winn J. 2010. A bayesian framework to  
534 account for complex non-genetic factors in gene expression levels greatly  
535 increases power in eqtl studies. *PLoS Comput Biol*, **6**: e1000770.  
536 Listgarten J, Kadie C, Schadt EE, Heckerman D. 2010. Correction for  
537 hidden confounders in the genetic analysis of gene expression. *Proc Natl  
538 Acad Sci U S A*, **107**: 16465–16470.  
539 Kang HM, Ye C, Eskin E. 2008. Accurate discovery of expression  
540 quantitative trait loci under confounding from spurious and genuine  
541 regulatory hotspots. *Genetics*, **180**: 1909–1925.  
542 Mostafavi S, Battle A, Zhu X, Urban AE, Levinson D, Montgomery SB,  
543 Koller D. 2013. Normalizing rna-sequencing data by modeling hidden  
544 covariates with prior knowledge. *PLoS One*, **8**: e68141.  
545 Battle A, *et al.* 2014. Characterizing the genetic basis of transcriptome  
546 diversity through rna-sequencing of 922 individuals. *Genome Res*, **24**:  
547 14–24.  
548 Alter O, *et al.* 2000. Singular value decomposition for genome-wide  
549 expression data processing and modeling. *Proceedings of the National  
550 Academy of Sciences*, **97**: 10101–10106.  
551 Alter O, Brown PO, Botstein D. 2003. Generalized singular value  
552 decomposition for comparative analysis of genome-scale expression data  
553 sets of two different organisms. *Proceedings of the National Academy of  
554 Sciences*, **100**: 3351–3356.  
555 Mostafavi S, *et al.* 2013. Normalizing rna-sequencing data by modeling  
556 hidden covariates with prior knowledge. *PLoS One*, **8**: e68141.  
557 Friedman JH. 1987. Exploratory projection pursuit. *Journal of the  
558 American statistical association*, **82**: 249–266.  
559 Bergstra J, Bengio Y. 2012. Random search for hyper-parameter  
560 optimization. *Journal of Machine Learning Research*, **13**: 281–305.  
561 Snoek J, Larochelle H, Adams RP. 2012. Practical bayesian optimization of  
562 machine learning algorithms. In *Advances in neural information  
563 processing systems*, S. 2951–2959.  
564 Basu K, Ghosh S. 2017. Analysis of thompson sampling for gaussian process  
565 optimization in the bandit setting. *arXiv preprint arXiv:1705.06808*.  
566 Agrawal S, Goyal N. 2013. Thompson sampling for contextual bandits with  
567 linear payoffs. In *International Conference on Machine Learning*, S.  
568 127–135.  
569 Hernández-Lobato JM, Hoffman MW, Ghahramani Z. 2014. Predictive  
570 entropy search for efficient global optimization of black-box functions.  
571 In *Advances in neural information processing systems*, S. 918–926.  
572 Agrawal S, Goyal N. 2013. Further optimal regret bounds for thompson  
573 sampling. In *Artificial Intelligence and Statistics*, S. 99–107.  
574 Rasmussen CE. 2004. Gaussian processes in machine learning. In *Advanced  
575 lectures on machine learning*, S. 63–71. Springer.  
576 Rahimi A, Recht B. 2008. Random features for large-scale kernel machines.  
577 In *Advances in neural information processing systems*, S. 1177–1184.  
578 Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, Hasz R,  
579 Walters G, Garcia F, Young N, *et al.* 2013. The genotype-tissue  
580 expression (gtex) project. *Nature genetics*, **45**: 580–585.  
581 Subramanian A, Tamayo P, *et al.* 2005. Gene set enrichment analysis: a  
582 knowledge-based approach for interpreting genome-wide expression  
583 profiles. *Proc Natl Acad Sci U S A*, **102**: 15545–15550.  
584 Williams A, Spilianakis CG, Flavell RA. 2010. Interchromosomal association  
585 and gene regulation in trans. *Trends in genetics*, **26**: 188–197.  
586 Petros Z, Lee MTM, Takahashi A, Zhang Y, Yimer G, Habtewold A,  
587 Schuppe-Koistinen I, Mushiroda T, Makonnen E, Kubo M, *et al.* 2017.  
588 Genome-wide association and replication study of hepatotoxicity induced  
589 by antiretrovirals alone or with concomitant anti-tuberculosis drugs.  
590 *OmicS: a journal of integrative biology*, **21**: 207–216.  
591 de Lange KM, Moutsianas L, Lee JC, Lamb CA, Luo Y, Kennedy NA,  
592 Jostins L, Rice DL, Gutierrez-Achury J, Ji SG, *et al.* 2017.

Genome-wide association study implicates immune activation of multiple  
integrin genes in inflammatory bowel disease. *Nature genetics*, **49**: 256.  
Astele WJ, Elding H, Jiang T, Allen D, Ruklisa D, Mann AL, Mead D,  
Bouman H, Riveros-Mckay F, Kostadima MA, *et al.* 2016. The allelic  
landscape of human blood cell trait variation and links to common  
complex disease. *Cell*, **167**: 1415–1429.

593  
594  
595  
596  
597  
598